

FAST TO TRAIN, FAST TO SAMPLE: STABLE VELOCITY FOR FLOW MATCHING

Anonymous authors

Paper under double-blind review

ABSTRACT

We revisit flow matching from a variance-centric perspective. Although conditional flow matching (CFM) is theoretically elegant, its use of single-sample conditional velocities introduces high variance, which can destabilize optimization and slow convergence. We demonstrate that this behavior induces two distinct regimes: a high-variance regime that hinders training and a low-variance regime where conditional and true velocities are nearly identical, thereby enabling analytical sampling shortcuts. Motivated by these insights, we introduce the **Stable Velocity** framework to improve both the training and sampling processes of flow matching. For training, we propose *Stable Velocity Matching (StableVM)*, a variance-reduced objective that preserves CFM’s global optima while significantly improving stability and convergence in the high-variance regime. For sampling, we introduce *Stable Velocity Sampling (StableVS)*, a “free lunch” acceleration method that leverages the low-variance regime to achieve faster generation without requiring finetuning. Experiments on SiT-XL trained on ImageNet, as well as on several large pretrained models (SD3, SD3.5, Flux, and Wan2.2), show consistent improvements in training convergence and more than $2\times$ faster sampling while maintaining high fidelity.

1 INTRODUCTION

Recent years have seen significant progress in generative modeling with the advent of diffusion (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020; Karras et al., 2022b), flow matching (Lipman et al., 2022; Liu et al., 2022), and stochastic interpolants (Albergo et al., 2023; Ma et al., 2024; Yu et al., 2024). These approaches model the transformation of a simple prior distribution, e.g., Gaussian noise, into a complex data distribution by a time-dependent stochastic or deterministic process. Through this lens, diffusion and flow-based models can be unified as different instantiations of stochastic interpolants, providing a principled foundation for large-scale generative modeling. Their success has fueled breakthroughs across diverse applications such as high-fidelity image synthesis (Labs, 2024; Esser et al., 2024), image restoration (Rombach et al., 2022; Lin et al., 2024), and video generation (Brooks et al., 2024; Wan et al., 2025).

Among these advances, the Conditional Flow Matching (CFM) framework (Tong et al., 2023) has emerged as an elegant formulation. Instead of explicitly simulating the forward stochastic differential equation (SDE) or probability flow ordinary differential equation (PF-ODE), CFM directly trains a neural network to predict the conditional velocity field that transports intermediate states along the interpolant. At convergence, the learned neural velocity field provably coincides with the true probability flow (Tong et al., 2023; Vincent, 2011; Xu et al., 2023), offering both theoretical soundness and practical scalability. Consequently, CFM has become a cornerstone objective for training flow and diffusion-based generative models.

Despite its theoretical elegance, CFM has a critical, yet often overlooked, weakness: the variance of its training target. In practice, the conditional velocity $v_t(x_t | x_0)$ used in CFM is only a single-sample Monte Carlo estimate of the true velocity $v_t(x_t)$. This introduces potentially high variance, particularly at timesteps where the marginal distribution remains close to the simple prior. High-variance targets can destabilize

047 optimization, slow convergence, and create a mismatch between the empirical training dynamics and the
 048 ideal global optimum. While empirical studies have repeatedly hinted at such variance-induced inefficiencies
 049 in diffusion training (Karras et al., 2022a; Choi et al., 2022; Xu et al., 2023), a general variance-theoretic
 050 understanding within flow matching and broader stochastic interpolant frameworks has been lacking.

051 In this work, we develop a *variance-based perspective* on stochastic interpolants. By explicitly quantifying
 052 the variance of conditional velocities, we uncover a two-regime structure that governs training and inference:
 053 1) The *high-variance regime* arises at timesteps closer to the prior (high-noise), where conditional
 054 velocity deviates substantially from the true velocity field. 2) The *low-variance regime* arises at timesteps
 055 closer to the data distribution (low-noise), where the conditional velocity aligns closely with the true velocity
 056 field. Building on these observations, we propose **Stable Velocity**, a novel framework for both training and
 057 sampling in flow matching. In particular, for training, we propose **Stable Velocity Matching (StableVM)**,
 058 which provably shares the same global optima as CFM but reduces variance by aggregating multiple refer-
 059 ence samples. For sampling, our variance analysis reveals that once variance is negligible, the dynamics
 060 become stable and admit closed-form simplifications that reduce the number of steps required for infer-
 061 ence without sacrificing sample quality. We call this "free-lunch" acceleration **Stable Velocity Sampling**
 062 (**StableVS**). We validate our approach across SiT-XL trained on ImageNet, and state-of-the-art pretrained
 063 generative models. StableVM consistently improves training convergence (Fig. 3 and 5), while StableVS
 064 accelerates inference by more than $2\times$ (e.g., 50 steps to 25 steps) for the latest flow-based models, includ-
 065 ing Stable Diffusion 3, Stable Diffusion 3.5, Flux, and Wan2.2, without quality degradation. Our main
 066 contributions are summarized as follows:

- 067 • **Variance analysis of stochastic interpolants.** We provide the first systematic study of variance on
 068 stochastic interpolants, identifying a two-regime structure that governs training and inference.
- 069 • **Stable Velocity Matching (StableVM).** We propose an unbiased, variance-reduced training objective that
 070 provably retains the same global optima of CFM while dramatically lowering training variance.
- 071 • **Stable Velocity Sampling (StableVS).** We exploit the low-variance regime to derive simplified dynamics
 072 that enable accelerated and stable sampling in pretrained models without additional finetuning.

074 2 PRELIMINARIES

075 We present a brief overview of flow matching and stochastic interpolants. Their detailed connections to
 076 score-based diffusion models are provided in Appendix A.

077 **Flow Matching and Stochastic Interpolants.** Given samples from an unknown data distribution $q(\mathbf{x}_0)$ over
 078 \mathbb{R}^d , the goal of generative modeling is to learn a model capable of generating new samples from q . Flow
 079 matching (Lipman et al., 2022; Liu et al., 2022; Lipman et al., 2024b) and stochastic interpolants (Albergo
 080 et al., 2023; Albergo & Vanden-Eijnden, 2023) approach this task by defining a continuous-time stochastic
 081 process. Starting with a data point $\mathbf{x}_0 \sim q$ and Gaussian noise $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$, the interpolant is defined as

$$082 \mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \varepsilon, \quad \alpha_0 = \sigma_1 = 1, \alpha_1 = \sigma_0 = 0, \quad (1)$$

083 where α_t and σ_t are differentiable functions satisfying $\alpha_t^2 + \sigma_t^2 > 0$ for all $t \in [0, 1]$. In practice, most
 084 works (Ma et al., 2024; Yu et al., 2024; Lipman et al., 2022; Liu et al., 2022; Lipman et al., 2024b) adopt a
 085 simple linear interpolant $\alpha_t = 1-t$, $\sigma_t = t$. This process induces a conditional and a marginal velocity field,

$$086 \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) = \frac{\sigma_t'}{\sigma_t}(\mathbf{x}_t - \alpha_t \mathbf{x}_0) + \alpha_t' \mathbf{x}_0 \quad (2) \quad \mathbf{v}_t(\mathbf{x}_t) = \mathbb{E}_{p_t(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)] \quad (3)$$

087 such that the probability flow ordinary differential equation (PF-ODE)

$$088 d\mathbf{x}_t = \mathbf{v}_t(\mathbf{x}_t) dt \quad (4)$$

089 induces marginal distributions that match that of Eq. (1) for all time $t \in [0, 1]$. Sampling from the model
 090 thus reduces to solving the PF-ODE in Eq. (4) using standard ODE solvers, e.g., Euler integration, starting
 091 from Gaussian noise $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$ (Lipman et al., 2022; Ma et al., 2024).
 092
 093

In addition, there exists a reverse stochastic differential equation (SDE) whose marginal $p_t(x)$ coincides with that of the PF-ODE in Eq. (4), but with an added diffusion term (Ma et al., 2024):

$$d\mathbf{x}_t = \mathbf{v}_t(\mathbf{x}_t) dt - \frac{1}{2}w_t \mathbf{s}_t(\mathbf{x}_t) dt + \sqrt{w_t} d\bar{\mathbf{w}}_t, \quad (5)$$

where $\bar{\mathbf{w}}_t$ is a standard Wiener process in backward time, $\sqrt{w_t}$ is the diffusion coefficient, and the score $\mathbf{s}_t(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ can be re-expressed using the velocity field (Ma et al., 2024):

$$\mathbf{s}_t(\mathbf{x}_t) = \sigma_t^{-1}(\alpha_t \mathbf{v}_t(\mathbf{x}_t) - \alpha'_t \mathbf{x}_t) / (\alpha'_t \sigma_t - \alpha_t \sigma'_t). \quad (6)$$

Conditional Flow Matching (CFM). Training typically uses the CFM objective (Tong et al., 2023):

$$\min_{\theta} \mathbb{E}_{t, q(\mathbf{x}_0), p_t(\mathbf{x}_t | \mathbf{x}_0)} \lambda_t \|\mathbf{v}_{\theta}(\mathbf{x}_t, t) - \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)\|^2, \quad (7)$$

where λ_t is a positive weighting function, and $\mathbf{v}_{\theta}(\cdot, \cdot) : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a neural velocity field parameterized by θ . Here, the conditional path distribution is $p_t(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$. The minimizer of Eq. (7) is provably the true marginal velocity field (Tong et al., 2023; Xu et al., 2023; Ma et al., 2024):

$$\mathbf{v}_{\theta}^*(\mathbf{x}_t, t) = \mathbb{E}_{p_t(\mathbf{x}_0 | \mathbf{x}_t)} [\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)] = \mathbf{v}_t(\mathbf{x}_t). \quad (8)$$

Variance of CFM. Although CFM provides an unbiased estimator of the true velocity field $\mathbf{v}_t(\mathbf{x}_t)$, it suffers from a key limitation: its training target, $\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)$, is only a single-sample Monte Carlo estimate of the expectation in Eq. (3). This introduces high variance (Owen, 2013; Elvira & Martino, 2021), which can slow convergence and degrade performance. Following Xu et al. (2023), we quantify this variance by the average trace of the conditional velocity variance at time t :

$$\begin{aligned} \mathcal{V}_{\text{CFM}}(t) &= \mathbb{E}_{p_t(\mathbf{x}_t)} [\text{Tr}(\text{Cov}_{p_t(\mathbf{x}_0 | \mathbf{x}_t)}(\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)))] \\ &= \mathbb{E}_{q(\mathbf{x}_0), p_t(\mathbf{x}_t | \mathbf{x}_0)} [\|\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2]. \end{aligned} \quad (9)$$

This definition mirrors the CFM loss in Eq. (7) without the expectation over t and with the neural velocity \mathbf{v}_{θ} replaced by the true velocity $\mathbf{v}_t(\mathbf{x}_t)$. Intuitively, a small $\mathcal{V}_{\text{CFM}}(t)$ means that $\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)$ is close to the true velocity $\mathbf{v}_t(\mathbf{x}_t)$, whereas a large value indicates strong deviations.

To better understand the behavior of $\mathcal{V}_{\text{CFM}}(t)$, we evaluate it on Gaussian Mixture Models (GMMs) and CIFAR-10 (Krizhevsky, 2009). As shown in Fig. 1, we make two consistent observations:

1. For all cases, the variance remains close to zero at small t but increases rapidly as t grows. This naturally divides the process into two regimes separated by a split point ξ : a *low-variance regime* ($0 \leq t < \xi$) with $\mathcal{V}_{\text{CFM}}(t) \approx 0$, and a *high-variance regime* ($\xi \leq t \leq 1$) where $\mathcal{V}_{\text{CFM}}(t)$ is significantly larger.
2. As the data dimensionality grows, the split point ξ shifts closer to 1, enlarging *low-variance regime* while compressing *high-variance regime*. At the same time, the overall variance magnitude increases.

Fig. 2 further illustrates these phenomena. In the *low-variance regime*, \mathbf{x}_t is often influenced by a single sample, so the conditional velocity nearly coincides with the true velocity, keeping $\mathcal{V}_{\text{CFM}}(t)$ small. In the

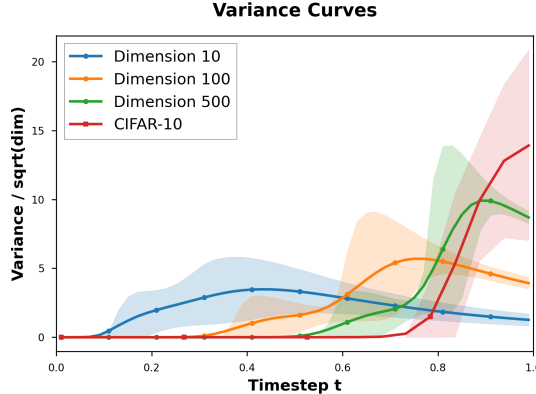


Figure 1: Variance curves of $\mathcal{V}_{\text{CFM}}(t)$ with 15%–85% quantile intervals, evaluated on GMMs with varying dimensionality and on CIFAR-10. The y -axis shows $\mathcal{V}_{\text{CFM}}(t)$ normalized by the square root of data dimension. More details are provided in Appendix E.1.4.

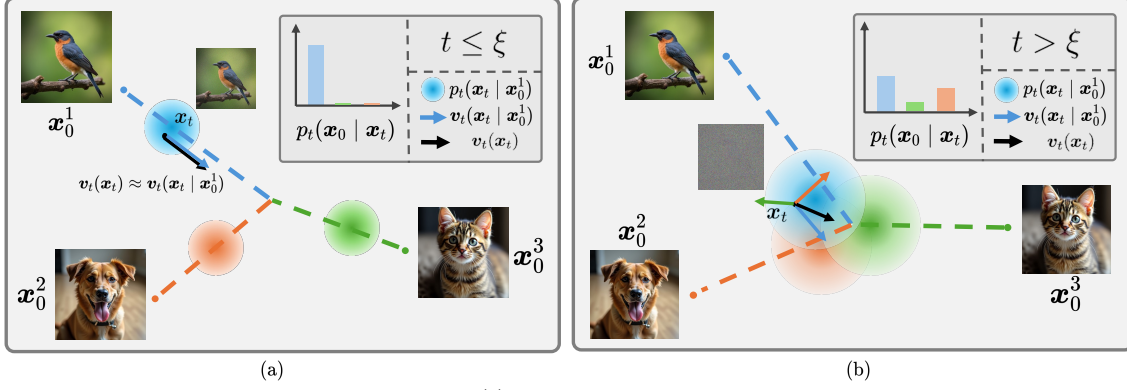


Figure 2: Illustration of CFM variance $\mathcal{V}_{\text{CFM}}(t)$. (a) The *low-variance regime* ($t \leq \xi$), where the posterior $p_t(\mathbf{x}_0 | \mathbf{x}_t)$ is sharply concentrated and the conditional velocity $\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)$ nearly coincides with the true velocity $\mathbf{v}_t(\mathbf{x}_t)$, yielding $\mathcal{V}_{\text{CFM}}(t) \approx 0$. (b) The *high-variance regime* ($t > \xi$), the posterior spreads over multiple reference samples, causing the conditional velocity to fluctuate and resulting in a large $\mathcal{V}_{\text{CFM}}(t)$.

high-variance regime, $\mathbf{x}_t \sim p_t$ is influenced by multiple samples, causing the conditional velocity to deviate substantially from the true velocity and leading to higher variances. As dimensionality increases, the distances between data points grow, delaying this mixing effect and shifting the split point ξ closer to 1.

These observations naturally suggest two key research directions: 1) In the *high-variance regime*, can we design methods to reduce variance to improve training convergence while preserving the correct global minimizer? 2) In the *low-variance regime*, where the velocity field already aligns closely with the true velocity field, can we develop sampling strategies that speed up generation without compromising performance?

3 VARIANCE-DRIVEN OPTIMIZATION OF TRAINING AND SAMPLING

In this section, we address the two key questions raised in Sec. 2 by introducing **Stable Velocity Matching (StableVM)**, a new variance-reduced training objective for stochastic interpolants, and **Stable Velocity Sampling (StableVS)**, an accelerated sampling method that does not require finetuning.

3.1 STABLE VELOCITY MATCHING: REDUCING TRAINING VARIANCE

Stable Velocity Matching (StableVM). We first introduce our novel training objective—StableVM and then show that this objective shares exactly the same global minimizer as CFM in Eq. (8). Furthermore, we quantify how its trace-of-variance decays as the reference batch size grows. Inspired by Xu et al. (2023), we introduce n reference samples $\{\mathbf{x}_0^i\}_{i=1}^n$, drawn i.i.d. from the data distribution $q(\mathbf{x}_0)$. We then define a composite conditional probability path $p_t^{\text{GMM}}(\mathbf{x}_t | \{\mathbf{x}_0^i\}_{i=1}^n) := \sum_{i=1}^n \frac{1}{n} p_t(\mathbf{x}_t | \mathbf{x}_0^i)$, which is essentially a mixture of conditional probabilities associated with each reference sample. We can derive the corresponding posterior distribution as in the following proposition. The proof is in Appendix D.1.

Proposition 1. *The posterior $p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) = \frac{1}{n} \sum_{i=1}^n (p_t(\mathbf{x}_0^i | \mathbf{x}_t) \prod_{j \neq i} q(\mathbf{x}_0^j))$.*

Based on the reference samples and p_t^{GMM} , we formulate our proposed unbiased training objective within the general stochastic interpolant framework. This objective function, denoted as $\mathcal{L}_{\text{StableVM}}$, is defined as:

$$\begin{aligned} \mathcal{L}_{\text{StableVM}}(\boldsymbol{\theta}, t) &= \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n, \mathbf{x}_t \sim p_t^{\text{GMM}}(\cdot | \{\mathbf{x}_0^i\}_{i=1}^n)} \left\| \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 \\ &= \mathbb{E}_{\mathbf{x}_t \sim p_t, \{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot | \mathbf{x}_t)} \left\| \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2. \end{aligned} \quad (10)$$

Our StableVM target is compatible with general stochastic interpolant framework. Under the special case of VP diffusion, it closely resembles the STF objective Eq. (20), but differs in the construction of the composite conditional probability path p_t^{GMM} . A detailed discussion is provided in Appendix B. As shown in the following theorem, the use of p_t^{GMM} is the essential ingredient that guarantees the unbiasedness of our target. The detailed proof is provided in Appendix D.2.

Theorem 1. (a) *The StableVM target is unbiased. That is, for any \mathbf{x}_t , we have*

$$\mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot|\mathbf{x}_t)} \left[\frac{\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right] = \mathbf{v}_t(\mathbf{x}_t).$$

(b) *The global minimizer $\mathbf{v}^*(\mathbf{x}_t, t)$ of the StableVM objective $\mathcal{L}_{\text{StableVM}}$ is the true velocity field $\mathbf{v}_t(\mathbf{x}_t)$.*

Variance of StableVM. Having established its unbiasedness, we now analyze the variance of StableVM, highlighting its significantly reduced training target variance compared to CFM. To quantify this, we follow Eq. (9) and consider the average trace-of-variance for StableVM:

$$\begin{aligned} \mathcal{V}_{\text{StableVM}}(t) &= \mathbb{E}_{\mathbf{x}_t \sim p_t} \left[\text{Tr} \left(\text{Cov}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot|\mathbf{x}_t)} \left(\sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right) \right) \right] \\ &= \mathbb{E}_{\mathbf{x}_t \sim p_t, \{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot|\mathbf{x}_t)} \left\| \mathbf{v}_t(\mathbf{x}_t) - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 \end{aligned} \quad (11)$$

These equivalent expressions provide several ways to represent the variance of the StableVM training target, characterizing the mean squared error between the estimated target and the true velocity $\mathbf{v}_t(\mathbf{x}_t)$.

In the following theorems, we compare $\mathcal{V}_{\text{StableVM}}$ against \mathcal{V}_{CFM} . The proofs are provided in Appendix D.3. We first show that $\mathcal{V}_{\text{StableVM}}$ is always upper bounded by \mathcal{V}_{CFM} .

Theorem 2. *Fix $t \in [0, 1]$. We always have $\mathcal{V}_{\text{StableVM}}(t) \leq \mathcal{V}_{\text{CFM}}(t)$.*

We then provide a stronger variance bound in the large sample regime. In particular, the upper bound gets stricter with a factor of $O(1/n)$.

Theorem 3. *Fix $t \in [0, 1]$. Let \mathbf{v}_t be bounded. Let $\varepsilon \in (0, 1)$. Assume*

$$M := \int_{\{\mathbf{x} : p_t(\mathbf{x}) \leq \varepsilon\}} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot|\mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] d\mathbf{x}_t < \infty.$$

Then, for large enough n , we have

$$\mathcal{V}_{\text{StableVM}}(t) \lesssim \frac{1}{n-1} \left(\frac{1}{\varepsilon} \cdot \mathcal{V}_{\text{CFM}}(t) + M \right).$$

We note that the (informal) assumption of n being large enough comes from the limits in distribution of the random variables as $n \rightarrow \infty$ as in Lemmas 3 and 4. We also note that we can choose a suitable ε in practice such that the set $\{\mathbf{x} : p_t(\mathbf{x}) \leq \varepsilon\}$ is “small” (in the sense that it has Lebesgue measure close to 0), which would lead to M being close to 0.

Extension to Class-Conditional Generation with Classifier-Free Guidance. Thus far, we have discussed variance reduction in the unconditional setting. Extending StableVM to conditional generation poses additional challenges. The main difficulty lies in the sparsity of conditional distributions: if we directly adopt the on-the-fly sampling strategy from Alg. 1, only a few reference samples will match a given label or prompt, leading to an insufficiently small effective reference set. To address this, we introduce a memory bank mechanism tailored for class-conditional generation. The memory bank is pre-populated with reference samples from the entire training dataset prior to training. During training, this memory bank is updated using a First-In, First-Out (FIFO) policy. This design guarantees that each class maintains a sufficiently large and diverse pool of reference samples, thereby stabilizing the training objective. The full procedure is detailed in Alg. 2.

Algorithm 1 Stable Velocity Matching

Require: Training iteration T , initial model \mathbf{v}_θ , dataset \mathcal{D} , learning rate η , ema rate α .

- 1: **for** $iter = 1 \dots T$ **do**
- 2: Sample a batch $\mathcal{B} = \{\mathbf{x}_0^i\}$ from \mathcal{D} , $|\mathcal{B}| > 1$
- 3: Uniformly sample time $t \sim q_t(t)$ from $[0, 1]$
- 4: Sample perturbed batch $\{\mathbf{x}_t^j\}_{j=1}^M$ from $p_{\text{GMM}}(\mathbf{x}_t^j | \{\mathbf{x}_0^i\}) = \sum_{\mathbf{x}_0 \in \mathcal{B}} \frac{1}{|\mathcal{B}|} p_t(\mathbf{x}_t^j | \mathbf{x}_0)$
- 5: Calculate stable vector field for all \mathbf{x}_t^j : $\mathbf{v}_B(\mathbf{x}_t^j) = \sum_{\mathbf{x}_0 \in \mathcal{B}} \frac{p_t(\mathbf{x}_t^j | \mathbf{x}_0)}{\sum_{\mathbf{y}_0 \in \mathcal{B}} p_t(\mathbf{x}_t^j | \mathbf{y}_0)} \cdot \mathbf{v}_t(\mathbf{x}_t^j | \mathbf{x}_0)$
- 6: Calculate loss for \mathbf{v}_θ : $\mathcal{L}(\theta) = \frac{1}{M} \sum_{j=1}^M \lambda(t) \|\mathbf{v}_\theta(\mathbf{x}_t^j, t) - \mathbf{v}_B(\mathbf{x}_t^j)\|^2$
- 7: Update model \mathbf{v}_θ : $\theta = \theta - \eta \nabla \mathcal{L}(\theta)$
- 8: **end for**
- 9: **return** \mathbf{v}_θ

3.2 STABLE VELOCITY SAMPLING: ACCELERATING SAMPLING AT LOW-VARIANCE REGIME

As discussed in Sec. 2 and shown in Fig 2, the variance $\mathcal{V}_{\text{CFM}}(t)$ is nearly zero in the *low-variance regime*. In this case, the true velocity $\mathbf{v}_t(\mathbf{x}_t)$ is effectively determined by a single data point \mathbf{x}_0 and can be well approximated by the conditional velocity field $\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)$. In other words, once the dominant data point \mathbf{x}_0 is identified from \mathbf{x}_t and $\mathbf{v}_t(\mathbf{x}_t)$, the analytical form of the velocity field—and hence its trajectory—becomes available. This observation enables stable and analytical simulation of the reverse dynamics between timesteps τ and t such that $\tau < t \leq \xi$. Specifically, we substitute the true velocity field $\mathbf{v}_t(\mathbf{x}_t)$ with its conditional counterpart $\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)$ in the PF-ODE (Eq. (4)) and the reverse SDE (Eq. (5)), and then perform the corresponding integration. If the trained model $\mathbf{v}_\theta(\mathbf{x}_t, t)$ perfectly recovers $\mathbf{v}_t(\mathbf{x}_t, t)$, this yields a principled way to accelerate sampling without sacrificing performance.

StableVS for SDE. For the reverse SDE (Eq. (5)), this yields the following DDIM-style (Song et al., 2021) posterior distribution:

$$p_\tau(\mathbf{x}_\tau | \mathbf{x}_t, \mathbf{v}_t(\mathbf{x}_t)) = \mathcal{N}(\boldsymbol{\mu}_{\tau|t}, \beta_t^2 \mathbf{I}), \quad (12)$$

where $\beta_t = f_\beta \sigma_\tau$, the factor $f_\beta \in [0, 1]$, and the posterior mean is

$$\boldsymbol{\mu}_{\tau|t} = \left(\sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}} - \left(\alpha_\tau - \alpha_t \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}} \right) \cdot \frac{\sigma'_t / \sigma_t}{\alpha'_t - \alpha_t \sigma'_t / \sigma_t} \right) \mathbf{x}_t + \left(\alpha_\tau - \alpha_t \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}} \right) \cdot \frac{\mathbf{v}_t(\mathbf{x}_t)}{\alpha'_t - \alpha_t \sigma'_t / \sigma_t}, \quad (13)$$

where α'_t and σ'_t represent the time derivatives of α_t and σ_t . The derivation is provided in Appendix D.4.

StableVS for ODE. For the probability flow ODE, the exact solution at timestep τ is

$$\mathbf{x}_\tau = \sigma_\tau \left[\left(\frac{1}{\sigma_t} - \frac{\sigma'_t}{\sigma_t} \cdot \frac{I(t, \tau)}{C_t} \right) \mathbf{x}_t + \frac{I(t, \tau)}{C_t} \mathbf{v}_t(\mathbf{x}_t) \right], \quad (14)$$

where $C_t := \alpha'_t - \alpha_t \sigma'_t / \sigma_t$, $I(t, \tau) := \int_t^\tau C(s) / \sigma_s ds$. The derivation is in Appendix D.5. In the special case of linear interpolant (*i.e.*, $\alpha_t = 1 - t$, $\sigma_t = t$), setting $\beta_t = 0$ in Eq. (12) makes two samplers coincide:

$$\mathbf{x}_\tau = \mathbf{x}_t + (\tau - t) \mathbf{v}_t(\mathbf{x}_t). \quad (15)$$

For this interpolant, Eq. (15) shows that the PF-ODE trajectory reduces to a straight line with constant velocity $\mathbf{v}_t(\mathbf{x}_t)$, allowing exact integration via Euler steps of arbitrary size within the *low-variance regime*.

Table 1: ImageNet 256×256 evaluation with different bank sizes under StableVM. All metrics are measured with the SDE Euler-Maruyama sampler (NFE=250) with classifier-free guidance. \downarrow and \uparrow indicate whether lower or higher values are better, respectively.

Iterations	Model	FID \downarrow	sFID \downarrow	IS \uparrow	Precision \uparrow	Recall \uparrow
250k	SiT-XL/2	6.71	4.76	147.8	0.78	0.53
	+ StableVM (bank=256)	6.50	4.74	154.7	0.78	0.53
	+ StableVM (bank=512)	6.35	4.61	152.8	0.79	0.53
	+ StableVM (bank=1024)	5.85	4.57	157.4	0.79	0.53
500k	SiT-XL/2	4.01	4.53	188.1	0.80	0.55
	+ StableVM (bank=256)	3.66	4.52	197.0	0.81	0.55
	+ StableVM (bank=512)	3.39	4.46	204.1	0.81	0.56
	+ StableVM (bank=1024)	3.20	4.46	208.9	0.81	0.56
750k	SiT-XL/2	2.84	4.42	219.7	0.81	0.57
	+ StableVM (bank=256)	2.79	4.46	220.9	0.81	0.56
	+ StableVM (bank=512)	2.75	4.46	225.0	0.82	0.57
	+ StableVM (bank=1024)	2.50	4.53	238.1	0.82	0.56
1M	SiT-XL/2	2.41	4.43	236.2	0.82	0.58
	+ StableVM (bank=256)	2.38	4.57	242.4	0.82	0.57
	+ StableVM (bank=512)	2.39	4.41	238.5	0.82	0.57
	+ StableVM (bank=1024)	2.39	4.46	237.1	0.82	0.58

4 EXPERIMENTS

In this section, we empirically validate our Stable Velocity framework through extensive experiments. In Section 4.1, we demonstrate the effectiveness of our StableVM target by training an SiT-XL (Ma et al., 2024) on ImageNet (Deng et al., 2009). We also conduct a two-stage training and sampling experiment that exploits the high and low variance regimes. In Section 4.2, we verify our StableVS sampling acceleration on large-scale pretrained text-to-image models. Further details on experiments are provided in Appendix E.

4.1 EVALUATION ON STABLE VELOCITY MATCHING

To further assess the effectiveness of StableVM, we adopt SiT-XL/2 (Ma et al., 2024) as the backbone model and train the model on ImageNet (resized to 256×256) (Deng et al., 2009). We train the models for 1M iterations with a batch size of 256. The baseline is the standard CFM, while our approach replaces the training objective with StableVM under different per-class bank capacities K . Following prior work (Ma et al., 2024; Yu et al., 2024), we employ the Euler-Maruyama SDE sampler with $w_t = \sigma_t$ and fix the number of function evaluations (NFEs) to 250 for all methods to ensure a fair comparison. We also apply classifier-free guidance with a scale of $w = 1.35$ and interval-based scheduling (Kynkäänniemi et al., 2024). For evaluation, we report FID (Heusel et al., 2017), IS (Salimans et al., 2016), sFID (Nash et al., 2021), precision, and recall (Kynkäänniemi et al., 2019).

From Table 1, we see that StableVM consistently outperforms the baseline CFM across different training stages, yielding a lower FID and higher IS while maintaining comparable precision and recall. The improvements are particularly evident especially at earlier and mid-training stages (250k–750k iterations), where larger bank capacities (e.g., $K = 1024$) provide stronger gains. This indicates that a richer memory bank offers more stable and representative velocity targets, thereby accelerating and improving convergence.

Two-Stage Training and Sampling. To further validate the effectiveness of StableVM in *high-variance regime*, we propose to focus on only training t in *high-variance regime*, and use a pretrained model for sampling during *low-variance regime*. This two-stage design is also used in Wan2.2 (Wan et al., 2025), highlighting its scalability and potential for performance gains. In this setting, our StableVM objective naturally serves as an ideal choice for variance reduction. Concretely, we sample $t \in [\xi, 1]$ for training and compare the CFM loss with our StableVM loss (bank capacity $K = 256$).

At inference, the model trained on the high-variance regime is used for $t \in [\xi, 1]$, while a pre-trained SiT-XL/2 model is used for $t \in [0, \xi]$. Fig. 3 reports FID versus training steps, showing consistent improvements with StableVM.

Additional details and results, including experimental setup, ablation studies, and unconditional generation results on synthetic GMMs, are provided in Appendix E.1.

4.2 EVALUATION ON STABLE VELOCITY SAMPLING

We next evaluate the effectiveness of our proposed StableVS on large-scale pretrained text-to-image models. Fig. 4 presents qualitative comparisons using SD3.5 (Esser et al., 2024), where we show outputs from the default sampler with 50 and 25 steps, as well as our 25-step sampler under the same random seed. Compared to the default 25-step outputs, our method produces generations that remain visually closer to the 50-step references, preserving fine-grained details and semantic consistency even under a reduced step budget. Table 2 further quantifies this observation on *GenEval* (Ghosh et al., 2023) at 1024×1024 resolution across multiple pretrained models, including Flux-dev (Labs, 2024), SD3.5-Large and SD3-Medium (Esser et al., 2024). On non-reference metrics, our 25-step sampler consistently achieves performance closer to or surpassing the 50-step results compared to the default 25-step sampler. On reference metrics, StableVS provides dramatic improvements: PSNR nearly doubles and SSIM approaches the 50-step quality, confirming that our outputs remain perceptually closer to the ground-truth 50-step samples. Similar trends hold for Flux-dev and SD3-Medium, and these results validate that we can cut off redundant sampling steps in the *low-variance regime*. Additional experimental details and results are provided in Appendix E.2.

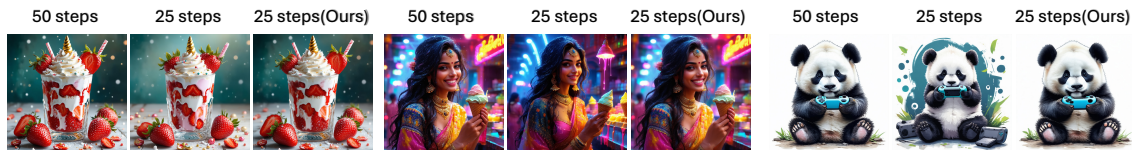


Figure 4: Visual comparisons of different prompts with SD3.5 (Esser et al., 2024). We show results using the default sampler with 50 and 25 steps, alongside our 25-step sampler under the same random seed. Compared to the default 25-step sampler, our method more closely aligns with the 50-step outputs. Zoom in for finer details. Additional qualitative comparisons are provided in Appendix G.

5 RELATED WORKS

Training Acceleration of Diffusion/Flow Models. Diffusion and flow matching models exhibit remarkable generative capabilities but incur substantial computational costs during training on high-resolution data. To mitigate this burden, prior work has pursued three complementary strategies: First, dimensionality reduction techniques aim to streamline training by compressing high-resolution data into lower-dimensional representations. Rombach et al. (2022) pioneered latent-space diffusion training, later enhanced by Chen et al. (2025) through improved compression ratios. Parallel efforts, such as the patch-wise score matching framework by Wang et al. (2023), further localized computational savings. A second direction integrates regularization with self-supervised learning objectives, as explored in Zheng et al. (2023); Yu et al. (2024), where auxiliary tasks stabilize diffusion model optimization. The third strategy focuses on improving the loss formulation. Studies like Karras et al. (2022b); Choi et al. (2022) address temporal inconsistencies in



Figure 3: Plot of FID for the two-stage experiments. Each model is trained up to 500k steps, and the checkpoints are evaluated every 40k steps. The data used for the plots can be found in Table 3.

Table 2: Evaluation results of multiple pretrained models on *GenEval* at 1024×1024 resolution, comparing the original method and ours. Metrics include non-reference metrics (overall, single-object, two-object, counting, colors, position, and color attribution) and reference metrics (PSNR, SSIM, LPIPS). Arrows (\uparrow / \downarrow) indicate whether higher or lower values are better.

Model	Non-reference metrics							Reference metrics		
	Overall \uparrow	Single \uparrow	Two \uparrow	Counting \uparrow	Colors \uparrow	Position \uparrow	Color Attr \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
<i>SD3.5-Large</i> (Esser et al., 2024)										
50steps	0.724	0.997	0.904	0.734	0.819	0.298	0.593	–	–	–
25steps	0.711	0.997	0.886	0.675	0.832	0.295	0.583	15.576	0.735	0.370
25steps (Ours)	0.718	0.994	0.904	0.706	0.806	0.300	0.600	30.954	0.958	0.050
<i>Flux-dev</i> (Labs, 2024)										
50steps	0.662	0.991	0.796	0.719	0.782	0.218	0.468	–	–	–
25steps	0.663	0.991	0.816	0.703	0.782	0.208	0.478	18.406	0.794	0.281
25steps (Ours)	0.657	0.988	0.816	0.706	0.766	0.218	0.450	29.187	0.937	0.068
<i>SD3-medium</i> (Esser et al., 2024)										
50steps	0.719	1.000	0.881	0.631	0.859	0.328	0.618	–	–	–
25steps	0.703	0.994	0.864	0.578	0.870	0.318	0.593	14.872	0.416	0.680
25steps (Ours)	0.713	0.997	0.861	0.625	0.854	0.315	0.625	28.493	0.940	0.077

training targets by proposing log-normal timestep sampling and signal-to-noise ratio (SNR)-aware weighting. Subsequent work Xu et al. (2023); Niedoba et al. (2024) employs self-normalized importance sampling (SNIS) estimators (Hesterberg, 1995) to reduce the variance of score estimation at the expense of bias. In this work, we extend the latter paradigm to the stochastic interpolants framework, introduce conditional mixture probabilities to remove bias, and demonstrate its efficacy in class-conditioned generation.

Sampling Acceleration of Diffusion/Flow Models. Reducing the number of sampling steps has become a key focus for both practical efficiency and theoretical analysis. Existing acceleration strategies can be broadly categorized into *training-required* and *training-free* approaches. Training-required methods typically aim to leverage additional learning to reduce steps. One approach is to distill pretrained multi-step diffusion models into few-step models (Salimans & Ho, 2022; Sauer et al., 2024). Other methods, such as consistency models (Song et al., 2023) and inductive moment matching (Zhou et al., 2025), enforce self-consistency of network outputs across different time steps, obviating the need for explicit distillation. Mean flows (Geng et al., 2025) take a related approach by predicting the average velocity along a trajectory, enabling one-step sampling. Training-free approaches, on the other hand, focus on designing more accurate numerical solvers for the reverse ODE. DDIM (Song et al., 2021) formulates sampling as a deterministic ODE and applies first-order integration, while high-order solvers such as DPM-Solver (Lu et al., 2022) and DPM-Solver++ (Lu et al., 2025) exploit the semi-linear structure to achieve higher accuracy in few-step generation. UniPC (Zhao et al., 2023) further unifies predictor–corrector schemes for diverse noise schedules, significantly reducing the quality gap relative to full-trajectory sampling. Our stable velocity sampling method is training-free, using analytical simplification in *low-variance regime* to reduce steps and remains compatible with other training-free solvers by employing them in *high-variance regime*.

6 CONCLUSION

In this work, we present a variance-based perspective on stochastic interpolants and uncover a two-regime structure that impacts both training and sampling. In the *high-variance regime*, conditional velocity estimates exhibit large variances, causing unstable and inefficient optimization. In contrast, in the *low-variance regime*, the estimator closely aligns with the true velocity, enabling analytical simplifications that accelerate sampling. Building on these insights, we introduce the **Stable Velocity** framework, which comprises: i) **Stable Velocity Matching (StableVM)**, a variance-reduced training objective that stabilizes optimization in the high-variance regime; ii) **Stable Velocity Sampling (StableVS)**, which leverages the low-variance regime to achieve faster, high-fidelity generation. Extensive evaluations on standard benchmarks and diverse pretrained models show that our approach improves both training convergence and sampling efficiency.

REFERENCES

- Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *11th International Conference on Learning Representations, ICLR 2023, 2023*.
- Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A ViT backbone for diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>.
- Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models, 2025. URL <https://arxiv.org/abs/2410.10733>.
- Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11472–11481, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *Advances in Neural Information Processing Systems*, 2021.
- Víctor Elvira and Luca Martino. Advances in importance sampling. *arXiv preprint arXiv:2102.05407*, 2021.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- Shanghai Gao, Pan Zhou, Ming-Ming Cheng, and Shuicheng Yan. MDTv2: Masked diffusion transformer is a strong image synthesizer. *arXiv preprint arXiv:2303.14389*, 2023.
- Zhengyang Geng, Mingyang Deng, Xingjian Bai, J Zico Kolter, and Kaiming He. Mean flows for one-step generative modeling. *arXiv preprint arXiv:2505.13447*, 2025.
- Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023.
- Ali Hatamizadeh, Jiamei Song, Guilin Liu, Jan Kautz, and Arash Vahdat. DiffiT: Diffusion vision transformers for image generation. In *European Conference on Computer Vision*, 2024.
- Tim Hesterberg. Weighted average importance sampling and defensive mixture distributions. *Technometrics*, 37(2):185–194, 1995.

- 470 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs
471 trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural*
472 *Information Processing Systems*, 2017.
- 473
474 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural*
475 *information processing systems*, 33:6840–6851, 2020.
- 476
477 Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans.
478 Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23
479 (47):1–33, 2022.
- 480
481 Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. Simple diffusion: End-to-end diffusion for high
482 resolution images. In *International Conference on Machine Learning*, 2023.
- 483
484 Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu,
485 Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and
486 Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. In *Proceedings of the*
IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024.
- 487
488 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based
489 generative models. In *Advances in Neural Information Processing Systems*, 2022a.
- 490
491 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based
492 generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022b.
- 493
494 Diederik Kingma and Ruiqi Gao. Understanding diffusion objectives as the ELBO with simple data aug-
495 mentation. *Advances in Neural Information Processing Systems*, 2024.
- 496
497 Diederik P Kingma. Adam: A method for stochastic optimization. In *International Conference on Learning*
Representations, 2015.
- 498
499 Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 2009.
- 500
501 Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision
502 and recall metric for assessing generative models. In *Advances in Neural Information Processing Systems*,
2019.
- 503
504 Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Ap-
505 plying guidance in a limited interval improves sample and distribution quality in diffusion models. *arXiv*
preprint arXiv:2404.07724, 2024.
- 506
507 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- 508
509 E.L. Lehmann and J.P. Romano. *Testing Statistical Hypotheses*. Springer Texts in Statistics Series. Springer
510 International Publishing, 2023. ISBN 9783030705800. URL [https://books.google.com.hk/
books?id=NIDzwwEACAAJ](https://books.google.com.hk/books?id=NIDzwwEACAAJ).
- 511
512 Xinqi Lin, Jingwen He, Ziyang Chen, Zhaoyang Lyu, Bo Dai, Fanghua Yu, Yu Qiao, Wanli Ouyang, and Chao
513 Dong. Diffbir: Toward blind image restoration with generative diffusion prior. In *European conference*
on computer vision, pp. 430–448. Springer, 2024.
- 514
515 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for
516 generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

- 517 Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky T. Q. Chen,
518 David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code, 2024a. URL <https://arxiv.org/abs/2412.06264>.
519
520
- 521 Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen,
522 David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint*
523 *arXiv:2412.06264*, 2024b.
- 524 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer
525 data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
526
- 527 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver
528 for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing*
529 *Systems*, 35:5775–5787, 2022.
- 530 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for
531 guided sampling of diffusion probabilistic models. *Machine Intelligence Research*, pp. 1–22, 2025.
532
- 533 Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie.
534 Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In
535 *European Conference on Computer Vision*, pp. 23–40. Springer, 2024.
- 536 Charlie Nash, Jacob Menick, Sander Dieleman, and Peter W Battaglia. Generating images with sparse
537 representations. In *International Conference on Machine Learning*, 2021.
538
- 539 Matthew Niedoba, Dylan Green, Saeid Naderiparizi, Vasileios Lioutas, Jonathan Wilder Lavington, Xiaox-
540 uan Liang, Yunpeng Liu, Ke Zhang, Setareh Dabiri, Adam Ścibior, et al. Nearest neighbour score esti-
541 mators for diffusion generative models. In *Proceedings of the 41st International Conference on Machine*
542 *Learning*, pp. 38117–38144, 2024.
- 543
- 544 Art B Owen. *Monte Carlo theory, methods and examples*. Stanford, 2013.
- 545 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *IEEE International*
546 *Conference on Computer Vision*, 2023.
547
- 548 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
549 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer*
550 *vision and pattern recognition*, pp. 10684–10695, 2022.
- 551
- 552 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint*
553 *arXiv:2202.00512*, 2022.
- 554 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved
555 techniques for training GANs. In *Advances in Neural Information Processing Systems*, 2016.
556
- 557 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation.
558 In *European Conference on Computer Vision*, pp. 87–103. Springer, 2024.
- 559 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning
560 using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, 2015.
561
- 562 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International*
563 *Conference on Learning Representations*, 2021.

- 564 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
565 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
566 *arXiv:2011.13456*, 2020.
- 567 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International*
568 *Conference on Machine Learning*, pp. 32211–32252. PMLR, 2023.
- 570 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the
571 Inception architecture for computer vision. In *IEEE Conference on Computer Vision and Pattern Recog-*
572 *niton*, 2016.
- 573 Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-Brooks,
574 Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch
575 optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- 577 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23
578 (7):1661–1674, 2011.
- 579 Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig
580 Davaadorj, Dhruv Nair, Sayak Paul, Steven Liu, William Berman, Yiyi Xu, and Thomas Wolf. Diffusers:
581 State-of-the-art diffusion models. URL <https://github.com/huggingface/diffusers>.
- 583 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwei Yu, Haiming
584 Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint*
585 *arXiv:2503.20314*, 2025.
- 586 Zhendong Wang, Yifan Jiang, Huangjie Zheng, Peihao Wang, Pengcheng He, Zhangyang Wang, Weizhu
587 Chen, Mingyuan Zhou, et al. Patch diffusion: Faster and more data-efficient training of diffusion models.
588 *Advances in neural information processing systems*, 36:72137–72154, 2023.
- 590 Yilun Xu, Shangyuan Tong, and Tommi Jaakkola. Stable target field for reduced variance score estimation
591 in diffusion models, 2023. URL <https://arxiv.org/abs/2302.00670>.
- 592 Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining
593 Xie. Representation alignment for generation: Training diffusion transformers is easier than you think.
594 *arXiv preprint arXiv:2410.06940*, 2024.
- 596 Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector
597 framework for fast sampling of diffusion models. *Advances in Neural Information Processing Systems*,
598 36:49842–49869, 2023.
- 599 Hongkai Zheng, Weili Nie, Arash Vahdat, and Anima Anandkumar. Fast training of diffusion models with
600 masked transformers. *arXiv preprint arXiv:2306.09305*, 2023.
- 602 Hongkai Zheng, Weili Nie, Arash Vahdat, and Anima Anandkumar. Fast training of diffusion models with
603 masked transformers. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- 604 Linqi Zhou, Stefano Ermon, and Jiaming Song. Inductive moment matching. *arXiv preprint*
605 *arXiv:2503.07565*, 2025.
- 606
607
608
609
610

A CONNECTIONS TO VARIANCE-PRESERVING DIFFUSION MODEL

In diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020), the forward process can be modeled as an Itô SDE:

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t) dt + g(t) d\mathbf{w}, \quad (16)$$

where \mathbf{w} is the standard Wiener process, $\mathbf{f}(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is vector-valued function called drift coefficient of \mathbf{x}_t and $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar function known as the diffusion coefficient of \mathbf{x}_t . It gradually transforms the data distribution q to a known prior as time goes from $t = 0$ to 1. With stochastic process defined in Eq. (16), Variance-Preserving (VP) diffusion model (Ho et al., 2020; Song et al., 2020; Karras et al., 2022a) implicitly define both α_t and σ_t in Eq. (1) with an equilibrium distribution as prior. VP diffusion commonly chooses $\alpha_t = \cos(\frac{\pi}{2}t)$, $\sigma_t = \sin(\frac{\pi}{2}t)$.

The diffusion models is trained to estimate the score of the marginal at time t , $\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t)$, via a neural network. Specifically, the training objective is a weighted sum of the denoising score matching (DSM) (Vincent, 2011):

$$\min_{\theta} \mathbb{E}_{t, q(\mathbf{x}_0), p_t(\mathbf{x}_t | \mathbf{x}_0)} \lambda_t \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \mathbf{s}_t(\mathbf{x}_t | \mathbf{x}_0)\|^2, \quad (17)$$

where λ_t is a positive weighting function, $\mathbf{s}_{\theta}(\cdot, \cdot) : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a time-dependent vector field parametrized as neural network with parameters θ , the conditional probability path $p_t(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$, and the conditional score function $\mathbf{s}_t(\mathbf{x}_t | \mathbf{x}_0) = \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{x}_0)$. The Eq. (17) shares a very similar form with CFM target in Eq. (7). Also, similar to Eq. (8), the objective in Eq. (17) admits a closed-form minimizer (Song et al., 2020; Xu et al., 2023):

$$\mathbf{s}_{\theta}^*(\mathbf{x}_t, t) = \mathbb{E}_{p_t(\mathbf{x}_0 | \mathbf{x}_t)} [\mathbf{s}_t(\mathbf{x}_t | \mathbf{x}_0)] = \mathbf{s}_t(\mathbf{x}_t) \quad (18)$$

Here, the marginal probability path $p_t(\mathbf{x}_t)$ is a mixture of conditional probability paths $p_t(\mathbf{x}_t | \mathbf{x}_0)$ that vary with data points \mathbf{x}_0 , that is,

$$p_t(\mathbf{x}_t) = \int p_t(\mathbf{x}_t | \mathbf{x}_0) q(\mathbf{x}_0) d\mathbf{x}_0. \quad (19)$$

B COMPARISON WITH STABLE TARGET FIELD

The standard training objective for diffusion models (Ho et al., 2020; Song et al., 2020; Karras et al., 2022b) is based on *denoising score matching* (DSM) (Vincent, 2011), also suffers from high variance.

To address this issue, Xu et al. (2023) proposed the *Stable Target Field* (STF), which stabilizes training by leveraging a reference batch $\mathcal{B} = \{\mathbf{x}_0^i\}_{i=1}^n \sim q(\mathbf{x}_0)$. The STF objective is defined as

$$\mathcal{L}_{\text{STF}}(\theta, t) = \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q(\mathbf{x}_0), p_t(\mathbf{x}_t | \mathbf{x}_0^i)} \left\| \mathbf{v}_{\theta}(\mathbf{x}_t, t) - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \mathbf{s}_t(\mathbf{x}_t | \mathbf{x}_0^k) \right\|^2. \quad (20)$$

Unlike DSM ($n = 1$), STF forms a weighted average of scores over the reference batch, with weights determined by the conditional likelihoods $p_t(\mathbf{x}_t | \mathbf{x}_0^k)$. This reduces the covariance of the target by a factor of n , thereby lowering variance. While STF introduces bias, the minimizer of \mathcal{L}_{STF} is given by

$$\mathbf{v}^*(\mathbf{x}_t, t) = \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t), \{\mathbf{x}_0^i\}_{i=2}^n \sim q(\mathbf{x}_0)} \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \mathbf{s}_t(\mathbf{x}_t | \mathbf{x}_0^k), \quad (21)$$

which deviates from the true score $\mathbf{s}_t(\mathbf{x}_t)$. However, as $n \rightarrow \infty$, this bias vanishes and the weighted estimator converges to the true score.

Our approach differs from STF in three important ways:

- 658 1. **General framework.** We extend the variance analysis and variance-reduction strategy to the flow match-
 659 ing as well as stochastic interpolant framework, which generalizes beyond VP diffusion and exhibits a
 660 distinct variance structure.
 661
 662 2. **Unbiased objective.** Instead of relying on a finite-sample weighted average, we propose a mixture of
 663 conditional probabilities that eliminates bias while still achieving variance reduction.
 664
 665 3. **Class-conditional extension.** While STF does not naturally extend to class-conditional settings, we
 666 design a tailored algorithm that maintains variance reduction under classifier-free guidance, improving
 667 both convergence and training efficiency.
 668

Algorithm 2 Stable Velocity Matching with Classifier-free Guidance

670 **Require:** Training iterations T , model \mathbf{v}_θ , dataset \mathcal{D} , learning rate η , batch size B , number of classes C , per-class bank
 671 capacity K , CFG dropout prob p_{cfg}

672 **Initialize memory bank** $\mathcal{M} = \{\mathcal{M}_c\}_{c=0}^C$ $\triangleright c=0, \dots, C-1$ are classes; $c=C$ is the *unconditional* bucket
 673 **for** $c = 0, \dots, C$ **do**

674 $\mathcal{M}_c \leftarrow$ Prefilled FIFO queue of capacity K

675 **end for**

676 **for** $iter = 1 \dots T$ **do**

677 Sample times $\{t_i\}_{i=1}^B \sim q_t([0, 1])$

678 Uniformly sample input labels $\{y^i\}_{i=1}^B$ from $0, \dots, C-1$

679 $y^i \leftarrow C$ with probability p_{cfg}

Stable Velocity Matching Update

680 **for** $i = 1 \dots B$ **do**

681 Sample perturbed samples $\mathbf{x}_{t_i}^i$ from

$$682 \quad p_{\text{GMM}}(\mathbf{x}_{t_i}^i | \mathcal{M}_{y^i}) = \frac{1}{|\mathcal{M}_{y^i}|} \sum_{\mathbf{x}_0^{\text{ref}} \in \mathcal{M}_{y^i}} p_t(\mathbf{x}_{t_i}^i | \mathbf{x}_0^{\text{ref}})$$

683 Compute stable field

$$684 \quad \mathbf{v}_{\mathcal{M}_{y^i}}(\mathbf{x}_{t_i}^i) = \sum_{\mathbf{x}_0^{\text{ref}} \in \mathcal{M}_{y^i}} \frac{p_t(\mathbf{x}_{t_i}^i | \mathbf{x}_0^{\text{ref}})}{\sum_{\mathbf{y}_0^{\text{ref}} \in \mathcal{M}_{y^i}} p_t(\mathbf{x}_{t_i}^i | \mathbf{y}_0^{\text{ref}})} \mathbf{v}_t(\mathbf{x}_{t_i}^i | \mathbf{x}_0^{\text{ref}})$$

685
 686
 687
 688
 689
 690 **end for**

691 Calculate loss for \mathbf{v}_θ :

$$692 \quad \mathcal{L}(\theta) = \frac{1}{B} \sum_{j=1}^B \lambda(t_j) \|\mathbf{v}_\theta(\mathbf{x}_{t_j}^j, t_j, y^j) - \mathbf{v}_{\mathcal{M}_{y^j}}(\mathbf{x}_{t_j}^j)\|^2$$

693 $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}(\theta)$

Online Memory Bank Update

694 Sample $\mathcal{B} = \{(\mathbf{x}_0^i, y_i)\}_{i=1}^B$ from \mathcal{D}

695 **for** $i = 1, \dots, B$ **do**

696 Push \mathbf{x}_0^i to \mathcal{M}_{y_i} (evict oldest if full)

697 Push \mathbf{x}_0^i to \mathcal{M}_C

\triangleright Unconditional bucket stores all samples

698 **end for**

699 **end for**

700 **return** \mathbf{v}_θ

C EXTENSION TO CLASS-CONDITIONAL GENERATION

We further extend StableVM to the class-conditional generation setting, with the complete procedure summarized in Algorithm 2.

D PROOFS

D.1 p_t^{GMM} POSTERIOR (PROPOSITION 1)

Proposition 1. *The posterior $p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) = \frac{1}{n} \sum_{i=1}^n (p_t(\mathbf{x}_0^i | \mathbf{x}_t) \prod_{j \neq i} q(\mathbf{x}_0^j))$.*

Proof. This is a simple application of the Bayes rule. Note that we have

$$\begin{aligned}
 p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) &= \frac{p_t^{\text{GMM}}(\mathbf{x}_t | \{\mathbf{x}_0^i\}_{i=1}^n) \cdot \prod_{i=1}^n q(\mathbf{x}_0^i)}{p_t(\mathbf{x}_t)} \\
 &= \frac{1}{p_t(\mathbf{x}_t)} \left(\sum_{i=1}^n \frac{1}{n} \cdot p_t(\mathbf{x}_t | \mathbf{x}_0^i) \right) \left(\prod_{i=1}^n q(\mathbf{x}_0^i) \right) \\
 &= \frac{1}{n} \sum_{i=1}^n \left(\frac{p_t(\mathbf{x}_t | \mathbf{x}_0^i)}{p_t(\mathbf{x}_t)} \cdot \prod_{j=1}^n q(\mathbf{x}_0^j) \right) \\
 &= \frac{1}{n} \sum_{i=1}^n \left(p_t(\mathbf{x}_0^i | \mathbf{x}_t) \prod_{j \neq i} q(\mathbf{x}_0^j) \right),
 \end{aligned} \tag{22}$$

as desired. □

D.2 PROOF OF UNBIASEDNESS (THEOREM 1)

Theorem 1. (a) *The StableVM target is unbiased. That is, for any \mathbf{x}_t , we have*

$$\mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot | \mathbf{x}_t)} \left[\frac{\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right] = \mathbf{v}_t(\mathbf{x}_t).$$

(b) *The global minimizer $\mathbf{v}^*(\mathbf{x}_t, t)$ of the StableVM objective $\mathcal{L}_{\text{StableVM}}$ is the true velocity field $\mathbf{v}_t(\mathbf{x}_t)$.*

752 *Proof.* (a) Using Eq. (22), we obtain
 753
 754

$$\begin{aligned}
 & \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot | \mathbf{x}_t)} \left[\frac{\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right] \\
 &= \int \frac{1}{n \cdot p_t(\mathbf{x}_t)} \left(\sum_{i=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^i) \right) \left(\prod_{i=1}^n q(\mathbf{x}_0^i) \right) \cdot \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} d\mathbf{x}_0^{1:n} \\
 &= \frac{1}{n \cdot p_t(\mathbf{x}_t)} \int \left(\prod_{i=1}^n q(\mathbf{x}_0^i) \right) \left(\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) \right) d\mathbf{x}_0^{1:n} \\
 &= \frac{1}{n \cdot p_t(\mathbf{x}_t)} \sum_{k=1}^n \int \left(\prod_{i=1}^n q(\mathbf{x}_0^i) \right) p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^{1:n} \\
 &= \frac{1}{n \cdot p_t(\mathbf{x}_t)} \sum_{k=1}^n \int \left(\prod_{i \neq k} q(\mathbf{x}_0^i) \right) p_t(\mathbf{x}_t, \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^{1:n} \\
 &= \frac{1}{n \cdot p_t(\mathbf{x}_t)} \sum_{k=1}^n \left(\left(\prod_{i \neq k} \int q(\mathbf{x}_0^i) d\mathbf{x}_0^i \right) \cdot \int p_t(\mathbf{x}_t, \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^k \right) \\
 &= \frac{1}{n} \sum_{k=1}^n \int p_t(\mathbf{x}_0^k | \mathbf{x}_t) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^k = \mathbf{v}_t(\mathbf{x}_t),
 \end{aligned}$$

755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778 as desired.

779
 780 (b) Recall the StableVM objective Eq. (10)
 781

$$\begin{aligned}
 & \mathcal{L}_{\text{StableVM}}(\boldsymbol{\theta}, t) \\
 &= \mathbb{E}_{\mathbf{x}_t \sim p_t, \{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot | \mathbf{x}_t)} \left[\left\| \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \frac{\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{x}_t \sim p_t} \left[\int p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) \left\| \mathbf{v}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \frac{\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 d\mathbf{x}_0^{1:n} \right].
 \end{aligned}$$

782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793 For each \mathbf{x}_t , let
 794

$$\mathcal{L}_{\mathbf{x}_t}(\mathbf{v}) := \int p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) \left\| \mathbf{v} - \frac{\sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 d\mathbf{x}_0^{1:n}.$$

It suffices to show that for each \mathbf{x}_t , setting \mathbf{v} to $\mathbf{v}_t(\mathbf{x}_t)$ above minimizes $\mathcal{L}_{\mathbf{x}_t}$. Note that $\mathcal{L}_{\mathbf{x}_t}$ is differentiable and strictly convex in \mathbf{v} , so the minimizer \mathbf{v}^* must satisfy $\nabla \mathcal{L}_{\mathbf{x}_t}(\mathbf{v}^*) = 0$. It follows that

$$\begin{aligned}
0 &= 2 \int p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) \left(\mathbf{v}^* - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right) d\mathbf{x}_0^{1:n} \\
&= 2 \int p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) \mathbf{v}^* d\mathbf{x}_0^{1:n} \\
&\quad - 2 \int p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t) \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} d\mathbf{x}_0^{1:n} \\
&= 2\mathbf{v}^* - 2\mathbf{v}_t(\mathbf{x}_t),
\end{aligned}$$

where we have used part (a) in the last step. Therefore, $\mathbf{v}_t(\mathbf{x}_t)$ is the unique minimizer of $\mathcal{L}_{\mathbf{x}_t}$. This finishes the proof. \square

D.3 PROOFS OF THE VARIANCE BOUNDS (THEOREM 2 AND THEOREM 3)

In this section, we prove the variance reduction bound of our StableVM target as in Eq. (11). We first recall that the StableVM target is the following estimator

$$\hat{\mathbf{v}}_t := \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \in \mathbb{R}^d, \tag{23}$$

where $\mathbf{x}_t \sim p_t$ and $\{\mathbf{x}_0^i\}_{i=1}^n \sim p_t^{\text{GMM}}(\cdot | \mathbf{x}_t)$.

We first prove the weaker version of the variance bound.

Theorem 2. Fix $t \in [0, 1]$. We always have $\mathcal{V}_{\text{StableVM}}(t) \leq \mathcal{V}_{\text{CFM}}(t)$.

846 *Proof.* We have

$$\begin{aligned}
847 & \mathcal{V}_{\text{StableVM}}(t) \\
848 & = \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n} \mathbb{E}_{\mathbf{x}_t \sim p_t^{\text{GMM}}(\cdot | \{\mathbf{x}_0^i\}_{i=1}^n)} \left\| \mathbf{v}_t(\mathbf{x}_t) - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 \\
849 & = \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n} \left[\int \frac{1}{n} \left(\sum_{i=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^i) \right) \left\| \mathbf{v}_t(\mathbf{x}_t) - \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \right\|^2 d\mathbf{x}_t \right] \\
850 & \leq \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n} \left[\int \frac{1}{n} \left(\sum_{i=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^i) \right) \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} \left\| \mathbf{v}_t(\mathbf{x}_t) - \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) \right\|^2 d\mathbf{x}_t \right] \\
851 & = \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n} \left[\int \frac{1}{n} \sum_{k=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^k) \left\| \mathbf{v}_t(\mathbf{x}_t) - \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) \right\|^2 d\mathbf{x}_t \right] \\
852 & = \frac{1}{n} \cdot \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n} \left[\sum_{k=1}^n \int p_t(\mathbf{x}_t | \mathbf{x}_0^k) \left\| \mathbf{v}_t(\mathbf{x}_t) - \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) \right\|^2 d\mathbf{x}_t \right] \\
853 & = \frac{1}{n} \cdot \mathbb{E}_{\{\mathbf{x}_0^i\}_{i=1}^n \sim q^n} \left[\sum_{k=1}^n \mathbb{E}_{\mathbf{x}_t \sim p_t(\cdot | \mathbf{x}_0^k)} \left\| \mathbf{v}_t(\mathbf{x}_t) - \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) \right\|^2 \right] \\
854 & = \frac{1}{n} \sum_{k=1}^n \mathcal{V}_{\text{CFM}}(t) = \mathcal{V}_{\text{CFM}}(t),
\end{aligned}$$

855 where the inequality is by Jensen's inequality and the convexity of $\|\cdot\|^2$. □

856 Before proceeding to proving the stronger bound, we first give an alternative interpretation of the posterior $p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t)$. Consider the following procedure conditioned on \mathbf{x}_t :

- 857 1. Sample a latent index I uniformly from $\{1, \dots, n\}$.
- 858 2. Sample $\mathbf{x}_0^I \sim p_t(\cdot | \mathbf{x}_t)$ and $\mathbf{x}_0^j \sim q$ for all $j \neq I$.

859 We claim the following:

860 **Lemma 1.** (a) The joint distribution of $\{\mathbf{x}_0^i\}_{i=1}^n$ sampled from the above procedure conditioned on \mathbf{x}_t is exactly $p_t^{\text{GMM}}(\{\mathbf{x}_0^i\}_{i=1}^n | \mathbf{x}_t)$.

861 (b) $\{\mathbf{x}_0^i\}_{i=1}^n$ are independent conditioned on \mathbf{x}_t and I .

862 *Proof.* (a) Note that the joint distribution of $\{\mathbf{x}_0^i\}_{i=1}^n$ sampled from the above procedure is

$$\frac{1}{n} \sum_{i=1}^n \left(p_t(\mathbf{x}_0^i | \mathbf{x}_t) \prod_{j \neq i} q(\mathbf{x}_0^j) \right),$$

863 which matches the joint distribution given by the posterior of p_t^{GMM} .

864 (b) This is clear by construction. □

The following lemmas compute the variance of the StableVM estimator step-by-step.

Lemma 2. *We have*

$$\text{Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t) = \mathbb{E}[\text{Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t, I) | \mathbf{x}_t],$$

where the expectation on the right-hand side is over the random variable I .

Proof. By the law of total covariance, we have

$$\text{Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t) = \mathbb{E}_I[\text{Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t, I) | \mathbf{x}_t] + \text{Cov}_I(\mathbb{E}[\widehat{\mathbf{v}}_t | \mathbf{x}_t, I] | \mathbf{x}_t).$$

We claim that $\mathbb{E}[\widehat{\mathbf{v}}_t | \mathbf{x}_t, I]$ does not depend on I , so the second term above would be 0.

Let $i \in [n]$. We have

$$\mathbb{E}[\widehat{\mathbf{v}}_t | \mathbf{x}_t, I = i] = \int p_t(\mathbf{x}_0^i | \mathbf{x}_t) \prod_{j \neq i} q(\mathbf{x}_0^j) \cdot \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} d\mathbf{x}_0^{1:n}.$$

Let $\pi : \mathbb{R}^{nd} \rightarrow \mathbb{R}^{nd}$ that swaps \mathbf{x}_0^1 and \mathbf{x}_0^i , i.e. $\pi(\mathbf{x}_0^1, \dots, \mathbf{x}_0^i, \dots, \mathbf{x}_0^n) = (\mathbf{x}_0^i, \dots, \mathbf{x}_0^1, \dots, \mathbf{x}_0^n)$. Note that we have $|\det D\pi| = 1$ since the Jacobian $D\pi$ is a permutation matrix. Then by the change of variable formula, we have

$$\begin{aligned} \mathbb{E}[\widehat{\mathbf{v}}_t | \mathbf{x}_t, I = i] &= \int p_t(\mathbf{x}_0^1 | \mathbf{x}_t) \prod_{j \neq 1} q(\mathbf{x}_0^j) \cdot \sum_{k=1}^n \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{j=1}^n p_t(\mathbf{x}_t | \mathbf{x}_0^j)} d\mathbf{x}_0^{1:n} \\ &= \mathbb{E}[\widehat{\mathbf{v}}_t | \mathbf{x}_t, I = 1]. \end{aligned}$$

This shows that $\text{Cov}_I(\mathbb{E}[\widehat{\mathbf{v}}_t | \mathbf{x}_t, I] | \mathbf{x}_t) = 0$, which finishes the proof. \square

For the next two lemmas, we fix $\ell \in [d]$ and only consider the ℓ^{th} component $\widehat{\mathbf{v}}_t^{(\ell)}$ for simplicity. Define

$$V_{n-1} := \sum_{k \neq i} p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0^k), \quad P_{n-1} := \sum_{k \neq i} p_t(\mathbf{x}_t | \mathbf{x}_0^k),$$

$$v_i := p_t(\mathbf{x}_t | \mathbf{x}_0^i) \mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0^i), \quad p_i := p_t(\mathbf{x}_t | \mathbf{x}_0^i).$$

Then we have

$$\widehat{\mathbf{v}}_t^{(\ell)} = \frac{V_{n-1} + v_i}{P_{n-1} + p_i}.$$

Note that conditioned on \mathbf{x}_t and $I = i$, the random variables \mathbf{x}_0^k for $k \neq i$ are all independent and follow the data distribution q by Lemma 1.

Let us also first compute some expectations that will be useful later. We have for $k \neq i$,

$$\mathbb{E}[p_t(\mathbf{x}_t | \mathbf{x}_0^k) | \mathbf{x}_t, I = i] = \int p_t(\mathbf{x}_t | \mathbf{x}_0^k) q(\mathbf{x}_0^k) d\mathbf{x}_0^k = \int p_t(\mathbf{x}_t, \mathbf{x}_0^k) d\mathbf{x}_0^k = p_t(\mathbf{x}_t), \quad (24)$$

$$\begin{aligned} \mathbb{E}[p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) | \mathbf{x}_t, I = i] &= \int p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) q(\mathbf{x}_0^k) d\mathbf{x}_0^k \\ &= \int p_t(\mathbf{x}_t, \mathbf{x}_0^k) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^k \\ &= p_t(\mathbf{x}_t) \int \frac{p_t(\mathbf{x}_t, \mathbf{x}_0^k)}{p_t(\mathbf{x}_t)} \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^k \\ &= p_t(\mathbf{x}_t) \int p_t(\mathbf{x}_0^k | \mathbf{x}_t) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^k) d\mathbf{x}_0^k \\ &= p_t(\mathbf{x}_t) \mathbf{v}_t(\mathbf{x}_t), \end{aligned} \quad (25)$$

$$\mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} [\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) | \mathbf{x}_t] = \int p_t(\mathbf{x}_0 | \mathbf{x}_t) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) d\mathbf{x}_0 = \mathbf{v}_t(\mathbf{x}_t), \quad (26)$$

and

$$\mathbb{E}_{\mathbf{x}_0 \sim q} \left[\frac{p_t(\mathbf{x}_0 | \mathbf{x}_t)}{q(\mathbf{x}_0)} \right] = \int p_t(\mathbf{x}_0 | \mathbf{x}_t) d\mathbf{x}_0 = 1. \quad (27)$$

We then continue with the lemmas.

Lemma 3. *As $n \rightarrow \infty$, we have*

$$\begin{aligned} & \sqrt{n-1} \left(\frac{V_{n-1}}{P_{n-1}} - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right) \\ & \xrightarrow{d} \mathcal{N} \left(0, \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right] \right), \end{aligned}$$

where the random variables V_{n-1} and P_{n-1} are conditioned on \mathbf{x}_t and $I = i$.

Proof. Write

$$f(\mathbf{x}_0) := \mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0), \quad w(\mathbf{x}_0) := \frac{p_t(\mathbf{x}_0 | \mathbf{x}_t)}{q(\mathbf{x}_0)}.$$

Then we have

$$\begin{aligned} \frac{V_{n-1}}{P_{n-1}} &= \frac{\sum_{k \neq i} p_t(\mathbf{x}_t | \mathbf{x}_0^k) \mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0^k)}{\sum_{k \neq i} p_t(\mathbf{x}_t | \mathbf{x}_0^k)} = \frac{\sum_{k \neq i} p_t(\mathbf{x}_t | \mathbf{x}_0^k) f(\mathbf{x}_0^k)}{\sum_{k \neq i} p_t(\mathbf{x}_t | \mathbf{x}_0^k)} \\ &= \frac{\sum_{k \neq i} \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) q(\mathbf{x}_0^k)}{p_t(\mathbf{x}_t) q(\mathbf{x}_0^k)} f(\mathbf{x}_0^k)}{\sum_{k \neq i} \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^k) q(\mathbf{x}_0^k)}{p_t(\mathbf{x}_t) q(\mathbf{x}_0^k)}} = \frac{\sum_{k \neq i} \frac{p_t(\mathbf{x}_0^k | \mathbf{x}_t)}{q(\mathbf{x}_0^k)} f(\mathbf{x}_0^k)}{\sum_{k \neq i} \frac{p_t(\mathbf{x}_0^k | \mathbf{x}_t)}{q(\mathbf{x}_0^k)}} \\ &= \frac{\sum_{k \neq i} w(\mathbf{x}_0^k) f(\mathbf{x}_0^k)}{\sum_{k \neq i} w(\mathbf{x}_0^k)}. \end{aligned} \quad (28)$$

Recall from Lemma 1 that conditioned on \mathbf{x}_t and $I = i$, the \mathbf{x}_0^i 's are all independent. Therefore, Eq. (28) is a self-normalized importance sampling estimator (Chapter 9 of Owen (2013)) with importance distribution $q(\mathbf{x}_0)$, nominal distribution $p_t(\mathbf{x}_0 | \mathbf{x}_t)$, and importance weight ratio $w(\mathbf{x}_0)$.

Note that by Eq. (26) and Eq. (27), we have

$$\mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} [f(\mathbf{x}_0)] = \mathbf{v}_t^{(\ell)}(\mathbf{x}_t), \quad \mathbb{E}_{\mathbf{x}_0 \sim q} [w(\mathbf{x}_0)] = 1.$$

Following results using the delta method as in Lehmann & Romano (2023) (Theorem 11.2.14) and Owen (2013) (Eq. (9.8)), we have that

$$\sqrt{n-1} \left(\frac{V_{n-1}}{P_{n-1}} - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right) \xrightarrow{d} \mathcal{N} \left(0, \mathbb{E}_{\mathbf{x}_0 \sim q} \left[w(\mathbf{x}_0)^2 \left(f(\mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right] \right).$$

Expanding the variance term above gives us

$$\begin{aligned}
& \mathbb{E}_{\mathbf{x}_0 \sim q} \left[w(\mathbf{x}_0)^2 \left(f(\mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right] \\
&= \mathbb{E}_{\mathbf{x}_0 \sim q} \left[\frac{p_t(\mathbf{x}_0 | \mathbf{x}_t)^2}{q(\mathbf{x}_0)^2} \cdot \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right] \\
&= \int \frac{p_t(\mathbf{x}_0 | \mathbf{x}_t)^2}{q(\mathbf{x}_0)} \cdot \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 d\mathbf{x}_0 \\
&= \frac{1}{p_t(\mathbf{x}_t)} \int p_t(\mathbf{x}_0 | \mathbf{x}_t) \cdot \frac{p_t(\mathbf{x}_0 | \mathbf{x}_t) p_t(\mathbf{x}_t)}{q(\mathbf{x}_0)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 d\mathbf{x}_0 \\
&= \frac{1}{p_t(\mathbf{x}_t)} \int p_t(\mathbf{x}_0 | \mathbf{x}_t) \cdot p_t(\mathbf{x}_t | \mathbf{x}_0) \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 d\mathbf{x}_0 \\
&= \frac{1}{p_t(\mathbf{x}_t)} \cdot \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[p_t(\mathbf{x}_t | \mathbf{x}_0) \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right],
\end{aligned}$$

as desired. \square

Lemma 4. Assume v_i is bounded. Then we have

$$\text{Var} \left(\widehat{\mathbf{v}}_t^{(\ell)} | \mathbf{x}_t, I = i \right) \approx \frac{1}{n-1} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right]$$

for large n .

Proof. Let $g(x, y) = x/y$ so that $\widehat{\mathbf{v}}_t^{(\ell)} = g(V_{n-1} + v_i, P_{n-1} + p_i)$. Performing a Taylor expansion of g at the point (V_{n-1}, P_{n-1}) gives us

$$\begin{aligned}
\widehat{\mathbf{v}}_t^{(\ell)} &= g(V_{n-1}, P_{n-1}) + \nabla g(V_{n-1} + \lambda_{n-1} v_i, P_{n-1} + \lambda_{n-1} p_i) \begin{bmatrix} v_i \\ p_i \end{bmatrix} \\
&= \frac{V_{n-1}}{P_{n-1}} + \frac{v_i}{P_{n-1} + \lambda_{n-1} p_i} - \frac{V_{n-1} + \lambda_{n-1} v_i}{(P_{n-1} + \lambda_{n-1} p_i)^2} p_i
\end{aligned}$$

for some $[0, 1]$ -valued random variable λ_{n-1} . It follows that

$$\begin{aligned}
& \sqrt{n-1} \left(\widehat{\mathbf{v}}_t^{(\ell)} - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right) \\
&= \sqrt{n-1} \left(\frac{V_{n-1}}{P_{n-1}} - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right) + \frac{\sqrt{n-1} \cdot v_i}{P_{n-1} + \lambda_{n-1} p_i} - \sqrt{n-1} \cdot \frac{V_{n-1} + \lambda_{n-1} v_i}{(P_{n-1} + \lambda_{n-1} p_i)^2} \cdot p_i \quad (29)
\end{aligned}$$

For the second term in Eq. (29), we first note that by the law of large numbers and Eq. (24), we have $\frac{P_{n-1}}{n-1} \xrightarrow{p} p_t(\mathbf{x}_t)$ as $n \rightarrow \infty$. Also note that $\frac{\lambda_{n-1}}{n-1} \rightarrow 0$ as $n \rightarrow \infty$ deterministically. It follows that $\frac{1}{n-1} (P_{n-1} + \lambda_{n-1} p_i) \xrightarrow{p} p_t(\mathbf{x}_t)$. We also have $\frac{v_i}{\sqrt{n-1}} \rightarrow 0$ deterministically since v_t is bounded. Hence, we have

$$\frac{\sqrt{n-1} \cdot v_i}{P_{n-1} + \lambda_{n-1} p_i} = \frac{\frac{1}{\sqrt{n-1}} v_i}{\frac{1}{n-1} (P_{n-1} + \lambda_{n-1} p_i)} \xrightarrow{p} 0$$

by Slutsky's theorem, so the second term of Eq. (29) converges to 0 in probability.

For the third term in Eq. (29), we have $\frac{1}{(n-1)^2} (P_{n-1} + \lambda_{n-1} p_i)^2 \xrightarrow{P} p_t(\mathbf{x}_t)^2$ by Slutsky's theorem. At the same time, by the law of large numbers and Eq. (25), we have $\frac{V_{n-1}}{n-1} \xrightarrow{P} p_t(\mathbf{x}_t) \mathbf{v}_t^{(\ell)}(\mathbf{x}_t)$. Since \mathbf{v}_t is bounded, we have $\frac{\lambda_{n-1} v_i}{n-1} \rightarrow 0$ deterministically. We also have $\frac{p_i}{\sqrt{n-1}} \rightarrow 0$ deterministically. Therefore, we get that

$$\begin{aligned} \frac{1}{(n-1)^{3/2}} (V_{n-1} + \lambda_{n-1} v_i) p_i &= \frac{p_i}{\sqrt{n-1}} \cdot \left(\frac{V_{n-1}}{n-1} + \frac{\lambda_{n-1} v_i}{n-1} \right) \\ &\xrightarrow{P} 0 \cdot \left(p_t(\mathbf{x}_t) \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) + 0 \right) = 0, \end{aligned}$$

by Slutsky's theorem. Then by Slutsky's theorem again, we have

$$\sqrt{n-1} \cdot \frac{V_{n-1} + \lambda_{n-1} v_i}{(P_{n-1} + \lambda_{n-1} p_i)^2} \cdot p_i = \frac{\frac{1}{(n-1)^{3/2}} (V_{n-1} + \lambda_{n-1} v_i) p_i}{\frac{1}{(n-1)^2} (P_{n-1} + \lambda_{n-1} p_i)^2} \xrightarrow{P} 0,$$

so the third term in Eq. (29) also converges to 0 in probability.

Therefore, combining Lemma 3 and the results above using Slutsky's theorem, we conclude that

$$\sqrt{n-1} \left(\widehat{\mathbf{v}}_t^{(\ell)} - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right) \xrightarrow{d} \mathcal{N} \left(0, \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right] \right).$$

This means that $\widehat{\mathbf{v}}_t^{(\ell)}$ has an approximated variance

$$\frac{1}{n-1} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right],$$

as desired. \square

We are now ready to state and prove the main theorem.

Theorem 3. Fix $t \in [0, 1]$. Let \mathbf{v}_t be bounded. Let $\varepsilon \in (0, 1)$. Assume

$$M := \int_{\{\mathbf{x}: p_t(\mathbf{x}) \leq \varepsilon\}} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] d\mathbf{x}_t < \infty.$$

Then, for large enough n , we have

$$\mathcal{V}_{\text{StableVM}}(t) \lesssim \frac{1}{n-1} \left(\frac{1}{\varepsilon} \cdot \mathcal{V}_{\text{CFM}}(t) + M \right).$$

Proof. Recall from Eq. (11) that we have

$$\begin{aligned} \mathcal{V}_{\text{StableVM}}(t) &= \mathbb{E} \left[\text{Tr Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t) \right] \\ &= \mathbb{E} \left[\text{Tr} \left(\mathbb{E} \left[\text{Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t, I) | \mathbf{x}_t \right] \right) \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\text{Tr Cov}(\widehat{\mathbf{v}}_t | \mathbf{x}_t, I) | \mathbf{x}_t \right] \right], \end{aligned}$$

where the second line follows from Lemma 2, and the third line follows from the linearity of expectation. Note that the inner expectation is over the random variable I , and the outer expectation is over the random variable \mathbf{x}_t . From Lemma 4, we have

$$\text{Var} \left(\widehat{\mathbf{v}}_t^{(\ell)} | \mathbf{x}_t, I = i \right) \approx \frac{1}{n-1} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot | \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t | \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right].$$

1081 Then

$$\begin{aligned}
1082 & \text{Tr Cov}(\widehat{\mathbf{v}}_t \mid \mathbf{x}_t, I = i) = \sum_{\ell=1}^d \text{Var}(\widehat{\mathbf{v}}_t^{(\ell)} \mid \mathbf{x}_t, I = i) \\
1083 & \approx \frac{1}{n-1} \sum_{\ell=1}^d \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t \mid \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \left(\mathbf{v}_t^{(\ell)}(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t^{(\ell)}(\mathbf{x}_t) \right)^2 \right] \\
1084 & = \frac{1}{n-1} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t \mid \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right].
\end{aligned}$$

1085 Since i does not appear in the last line above, we get

$$\begin{aligned}
1086 & \mathcal{V}_{\text{StableVM}}(t) \approx \frac{1}{n-1} \mathbb{E}_{\mathbf{x}_t \sim p_t} \left[\mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t \mid \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] \right] \\
1087 & \leq \frac{1}{n-1} \mathbb{E}_{\mathbf{x}_t \sim p_t} \left[\frac{1}{p_t(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] \right],
\end{aligned}$$

1088 where we have used $p_t(\mathbf{x}_t \mid \mathbf{x}_0) \leq 1$ in the second line.

1089 Now we fix some $\varepsilon > 0$. Define $P_{>\varepsilon} := \{\mathbf{x} \in \mathbb{R}^d : p_t(\mathbf{x}) > \varepsilon\}$ and $P_{\leq\varepsilon} := \mathbb{R}^d \setminus P_{>\varepsilon}$. Then we have

$$\begin{aligned}
1090 & \mathcal{V}_{\text{StableVM}}(t) \lesssim \frac{1}{n-1} \int_{\mathbb{R}^d} \frac{\mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right]}{p_t(\mathbf{x}_t)} \cdot p_t(\mathbf{x}_t) d\mathbf{x}_t \\
1091 & = \frac{1}{n-1} \int_{P_{>\varepsilon}} \frac{\mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right]}{p_t(\mathbf{x}_t)} \cdot p_t(\mathbf{x}_t) d\mathbf{x}_t \\
1092 & \quad + \frac{1}{n-1} \int_{P_{\leq\varepsilon}} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] d\mathbf{x}_t \\
1093 & \leq \frac{1}{\varepsilon(n-1)} \int_{P_{>\varepsilon}} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] \cdot p_t(\mathbf{x}_t) d\mathbf{x}_t + \frac{M}{n-1} \\
1094 & \leq \frac{1}{\varepsilon(n-1)} \int_{\mathbb{R}^d} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] \cdot p_t(\mathbf{x}_t) d\mathbf{x}_t + \frac{M}{n-1} \\
1095 & = \frac{1}{\varepsilon(n-1)} \cdot \mathbb{E}_{\mathbf{x}_t \sim p_t, \mathbf{x}_0 \sim p_t(\cdot \mid \mathbf{x}_t)} \left[\|\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_0) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right] + \frac{M}{n-1} \\
1096 & = \frac{1}{n-1} \left(\frac{1}{\varepsilon} \cdot \mathcal{V}_{\text{CFM}}(t) + M \right),
\end{aligned}$$

1097 as desired. □

1100 D.4 SIMULATING THE REVERSE SDE IN LOW-VARIANCE REGIME

1101 [Ma et al. \(2024\)](#) show that the reverse-time SDE (Eq. (5)) with score function $\mathbf{s}_t(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ and
1102 arbitrary diffusion strength $w_t \geq 0$ yields the correct marginal density $p_t(\mathbf{x})$ at each time t . Furthermore, as
1103 established in [Anderson \(1982\)](#); [Ma et al. \(2024\)](#), if $\mathbf{x}_{t_1} \sim p_{t_1}(\mathbf{x})$, then the reverse-time solution \mathbf{x}_τ at any
1104 $\tau \in [0, t_1]$ is distributed according to the posterior:

$$1105 p_\tau(\mathbf{x}_\tau \mid \mathbf{x}_{t_1}) = \mathbb{E}_{p_{t_1}(\mathbf{x}_0 \mid \mathbf{x}_{t_1})} [p_\tau(\mathbf{x}_\tau \mid \mathbf{x}_0, \mathbf{x}_{t_1})] \approx p_\tau(\mathbf{x}_\tau \mid \mathbf{x}_0, \mathbf{x}_{t_1}). \quad (30)$$

Proposition 2. Let $\mathbf{x}_t \sim \mathcal{N}(\alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$ and $\mathbf{x}_\tau \sim \mathcal{N}(\alpha_\tau \mathbf{x}_0, \sigma_\tau^2 \mathbf{I})$, where $\tau < t$, and $\mathbf{x}_0 \sim p(\mathbf{x}_0)$ is the clean data sample. For any fixed variance parameter $\beta_t^2 \in (0, \sigma_\tau^2)$, define the posterior distribution as

$$p_\tau^{\alpha_t}(\mathbf{x}_\tau | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(k_t \mathbf{x}_t + \lambda_t \mathbf{x}_0, \beta_t^2 \mathbf{I}),$$

then the coefficients

$$k_t = \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}}, \quad \lambda_t = \alpha_\tau - \alpha_t \cdot \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}}$$

guarantee that the marginal of \mathbf{x}_τ is $\mathcal{N}(\alpha_\tau \mathbf{x}_0, \sigma_\tau^2 \mathbf{I})$.

Proof. We begin by expressing \mathbf{x}_t using the forward diffusion process:

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}).$$

We define the reverse model as a Gaussian conditional:

$$\mathbf{x}_\tau = k_t \mathbf{x}_t + \lambda_t \mathbf{x}_0 + \eta, \quad \eta \sim \mathcal{N}(0, \beta_t^2 \mathbf{I}).$$

Substituting \mathbf{x}_t yields:

$$\mathbf{x}_\tau = k_t(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon) + \lambda_t \mathbf{x}_0 + \eta = (k_t \alpha_t + \lambda_t) \mathbf{x}_0 + k_t \sigma_t \epsilon + \eta.$$

Hence, the conditional distribution of \mathbf{x}_τ given \mathbf{x}_0 is:

$$\mathbf{x}_\tau | \mathbf{x}_0 \sim \mathcal{N}((k_t \alpha_t + \lambda_t) \mathbf{x}_0, (k_t^2 \sigma_t^2 + \beta_t^2) \mathbf{I}).$$

To match the desired marginal $\mathbf{x}_\tau \sim \mathcal{N}(\alpha_\tau \mathbf{x}_0, \sigma_\tau^2 \mathbf{I})$, we require:

$$\begin{aligned} k_t \alpha_t + \lambda_t &= \alpha_\tau, \\ k_t^2 \sigma_t^2 + \beta_t^2 &= \sigma_\tau^2. \end{aligned}$$

Solving the second equation above for k_t , we obtain:

$$k_t = \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}}.$$

Substituting into first equation, we get:

$$\lambda_t = \alpha_\tau - \alpha_t \cdot \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}}.$$

Thus, the choice of k_t and λ_t ensures that the conditional distribution of \mathbf{x}_τ is consistent with the marginal. \square

Within this low variance area, we also have

$$\mathbf{v}_t(\mathbf{x}_t) \approx \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) = \frac{\sigma_t'}{\sigma_t} (\mathbf{x}_t - \alpha_t \mathbf{x}_0) + \alpha_t' \mathbf{x}_0, \quad (31)$$

Thus, given the velocity field $\mathbf{v}_t(\mathbf{x}_t)$ and the current state \mathbf{x}_t , the target \mathbf{x}_0 can be extracted as:

$$\mathbf{x}_0 = \frac{\mathbf{v}_t(\mathbf{x}_t) - \frac{\sigma_t'}{\sigma_t} \mathbf{x}_t}{\alpha_t' - \frac{\sigma_t'}{\sigma_t} \alpha_t} \quad (32)$$

1175 Plugging in this equation into the original expression, thus the posterior distribution with \mathbf{x}_0 eliminated via
 1176 $\mathbf{v}_t(\mathbf{x}_t)$, is given by:

$$1177 \quad p_\tau(\mathbf{x}_\tau | \mathbf{x}_t, \mathbf{v}_t(\mathbf{x}_t)) = \mathcal{N}(\boldsymbol{\mu}_{\tau|t}, \beta_t^2 \mathbf{I})$$

1178 where the posterior mean is explicitly:

$$1179 \quad \boldsymbol{\mu}_{\tau|t} = \left(\sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}} - \left(\alpha_\tau - \alpha_t \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}} \right) \cdot \frac{\sigma'_t}{\alpha'_t - \frac{\sigma'_t}{\sigma_t} \alpha_t} \right) \mathbf{x}_t + \left(\alpha_\tau - \alpha_t \sqrt{\frac{\sigma_\tau^2 - \beta_t^2}{\sigma_t^2}} \right) \cdot \frac{\mathbf{v}_t(\mathbf{x}_t)}{\alpha'_t - \frac{\sigma'_t}{\sigma_t} \alpha_t}$$

1180 Assuming $\alpha_t = 1 - t$ and $\sigma_t = t$, the DDIM-style posterior becomes:

$$1181 \quad p_\tau(\mathbf{x}_\tau | \mathbf{x}_t, \mathbf{v}_t(\mathbf{x}_t)) = \mathcal{N}(\boldsymbol{\mu}_{\tau|t}, \beta_t^2 \mathbf{I})$$

1182 with mean:

$$1183 \quad \boldsymbol{\mu}_{\tau|t} = \left(\sqrt{\frac{\tau^2 - \beta_t^2}{t^2}} + \left((1 - \tau) - (1 - t) \sqrt{\frac{\tau^2 - \beta_t^2}{t^2}} \right) \right) \mathbf{x}_t - \left((1 - \tau) - (1 - t) \sqrt{\frac{\tau^2 - \beta_t^2}{t^2}} \right) t \mathbf{v}_t(\mathbf{x}_t)$$

1184 If we set $\beta_t = 0$, we obtain the deterministic sampler:

$$1185 \quad \mathbf{x}_\tau = \mathbf{x}_t + (\tau - t) \mathbf{v}_t(\mathbf{x}_t)$$

1186 D.5 EXPLICIT PF-ODE SOLUTION IN LOW-VARIANCE REGIME

1187 In *low-variance regime* ($0 \leq t \leq \xi$), the conditional velocity field simplifies as $\mathbf{v}_t(\mathbf{x}_t) \approx \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)$.
 1188 We can thus derive explicit solutions to the Probability Flow ODE (PF-ODE) under both the stochastic
 1189 interpolant and VP diffusion frameworks. We consider the **Probability Flow ODE (PF-ODE)** under the
 1190 stochastic interpolant framework:

$$1191 \quad \frac{d\mathbf{x}_t}{dt} = \mathbf{v}_t(\mathbf{x}_t) \approx \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) = \frac{\sigma'_t}{\sigma_t} (\mathbf{x}_t - \alpha_t \mathbf{x}_0) + \alpha'_t \mathbf{x}_0, \quad (33)$$

1192 where α_t and σ_t define a stochastic interpolant, and \mathbf{x}_0 is the data point to be matched.

1193 **Closed-form of the Target \mathbf{x}_0** Given the velocity field $\mathbf{v}_t(\mathbf{x}_t)$ and the current state \mathbf{x}_t , the target \mathbf{x}_0 can
 1194 be extracted as:

$$1195 \quad \mathbf{x}_0 = \frac{\mathbf{v}_t(\mathbf{x}_t) - \frac{\sigma'_t}{\sigma_t} \mathbf{x}_t}{\alpha'_t - \frac{\sigma'_t}{\sigma_t} \alpha_t} = \frac{\mathbf{v}_t(\mathbf{x}_t) - \frac{\sigma'_t}{\sigma_t} \mathbf{x}_t}{C_t}, \quad (34)$$

1196 where we define the coefficient

$$1197 \quad C_t := \alpha'_t - \frac{\sigma'_t}{\sigma_t} \alpha_t.$$

1198 **Solving the PF-ODE from t_1 to $0 \leq \tau < t_1$** We aim to integrate the PF-ODE backward in time from a
 1199 known terminal state \mathbf{x}_{t_1} . The PF-ODE can be written as:

$$1200 \quad \frac{d\mathbf{x}_t}{dt} + a(t) \mathbf{x}_t = b(t), \quad \text{where } a(t) = -\frac{\sigma'_t}{\sigma_t}, \quad b(t) = C_t \mathbf{x}_0. \quad (35)$$

This is a linear nonhomogeneous first-order ODE. The integrating factor is:

$$\mu(t) = \exp\left(\int a(t)dt\right) = \exp\left(-\int \frac{\sigma'_t}{\sigma_t} dt\right) = \frac{1}{\sigma_t}.$$

Multiplying both sides by $\mu(t)$ yields:

$$\frac{d}{dt}\left(\frac{\mathbf{x}_t}{\sigma_t}\right) = \frac{C_t}{\sigma_t} \mathbf{x}_0.$$

Integrating both sides from t_1 to $\tau < t_1$:

$$\frac{\mathbf{x}_\tau}{\sigma_\tau} = \frac{\mathbf{x}_{t_1}}{\sigma_{t_1}} + \int_{t_1}^{\tau} \frac{C(s)}{\sigma_s} ds \cdot \mathbf{x}_0, \quad (36)$$

which gives:

$$\mathbf{x}_\tau = \sigma_\tau \left(\frac{\mathbf{x}_{t_1}}{\sigma_{t_1}} + I(t_1, \tau) \cdot \mathbf{x}_0 \right), \quad (37)$$

where

$$I(t_1, \tau) := \int_{t_1}^{\tau} \frac{C(s)}{\sigma_s} ds.$$

Substituting \mathbf{x}_0 in Closed Form We now substitute the expression for \mathbf{x}_0 evaluated at time t_1 :

$$\mathbf{x}_0 = \frac{\mathbf{v}_{t_1}(\mathbf{x}_{t_1}) - \frac{\sigma'_{t_1}}{\sigma_{t_1}} \mathbf{x}_{t_1}}{C_{t_1}}.$$

Substitute this into the solution:

$$\mathbf{x}_\tau = \sigma_\tau \left(\frac{\mathbf{x}_{t_1}}{\sigma_{t_1}} + I(t_1, \tau) \cdot \frac{\mathbf{v}_{t_1}(\mathbf{x}_{t_1}) - \frac{\sigma'_{t_1}}{\sigma_{t_1}} \mathbf{x}_{t_1}}{C_{t_1}} \right) \quad (38)$$

$$= \sigma_\tau \left[\left(\frac{1}{\sigma_{t_1}} - \frac{\sigma'_{t_1}}{\sigma_{t_1}} \cdot \frac{I(t_1, \tau)}{C_{t_1}} \right) \mathbf{x}_{t_1} + \frac{I(t_1, \tau)}{C_{t_1}} \mathbf{v}_{t_1}(\mathbf{x}_{t_1}) \right]. \quad (39)$$

Final Expression (Only in Terms of \mathbf{x}_{t_1})

$$\mathbf{x}_\tau = \sigma_\tau \left[\left(\frac{1}{\sigma_{t_1}} - \frac{\sigma'_{t_1}}{\sigma_{t_1}} \cdot \frac{I(t_1, \tau)}{C_{t_1}} \right) \mathbf{x}_{t_1} + \frac{I(t_1, \tau)}{C_{t_1}} \mathbf{v}_{t_1}(\mathbf{x}_{t_1}) \right] \quad (40)$$

This provides a fully explicit backward solution to the PF-ODE, depending only on \mathbf{x}_{t_1} and the velocity field $\mathbf{v}_{t_1}(\mathbf{x}_{t_1})$.

Special Case: Linear Interpolant For the linear interpolant with $\alpha_t = 1 - t$ and $\sigma_t = t$, we have:

$$\alpha'_t = -1, \quad \sigma'_t = 1, \quad \Rightarrow \quad C_t = -\frac{1}{t}, \quad \frac{C_t}{\sigma_t} = -\frac{1}{t^2}.$$

Then:

$$I(t_1, \tau) = \int_{t_1}^{\tau} -\frac{1}{s^2} ds = \frac{1}{\tau} - \frac{1}{t_1}.$$

Also note:

$$\frac{\sigma'_{t_1}}{\sigma_{t_1}} = \frac{1}{t_1}, \quad C_{t_1} = -\frac{1}{t_1}.$$

Plug into the general expression:

$$\mathbf{x}_\tau = \tau \left[\left(\frac{1}{t_1} - \frac{1}{t_1} \cdot \frac{1/\tau - 1/t_1}{-1/t_1} \right) \mathbf{x}_{t_1} + \left(\frac{1/\tau - 1/t_1}{-1/t_1} \right) \mathbf{v}_{t_1}(\mathbf{x}_{t_1}) \right] \quad (41)$$

$$= \tau \left[\left(\frac{1}{t_1} + \left(\frac{1}{\tau} - \frac{1}{t_1} \right) \right) \mathbf{x}_{t_1} + \left(\frac{1}{t_1} - \frac{1}{\tau} \right) \mathbf{v}_{t_1}(\mathbf{x}_{t_1}) \right] \quad (42)$$

$$= \mathbf{x}_{t_1} + (\tau - t_1) \mathbf{v}_{t_1}(\mathbf{x}_{t_1}). \quad (43)$$

Final Linear Interpolant Result

$$\boxed{\mathbf{x}_\tau = \mathbf{x}_{t_1} + (\tau - t_1) \mathbf{v}_{t_1}(\mathbf{x}_{t_1})} \quad (44)$$

In the case of the linear interpolant, the PF-ODE corresponds to a straight-line trajectory with constant velocity $\mathbf{v}_{t_1}(\mathbf{x}_{t_1})$, enabling exact integration via Euler steps of arbitrary size.

E EXPERIMENTAL DETAILS & RESULTS

E.1 EVALUATION ON STABLE VELOCITY MATCHING

E.1.1 EXPERIMENTAL SETUP

Training hyperparameters. Our implementation builds upon the official REPA codebase (Yu et al., 2024). We adopt SiT-XL (Ma et al., 2024) as the backbone architecture. The model is optimized using AdamW (Kingma, 2015) with a constant learning rate of 1×10^{-4} , $\beta_1 = 0.9$, $\beta_2 = 0.999$, and no weight decay. To accelerate training, we employ mixed-precision (FP16) arithmetic together with gradient clipping. Additionally, we use an adaptive EMA rate schedule following Lipman et al. (2024a). Specifically, at each training step, we increment a global counter `num_updates` and compute the EMA decay as

$$\text{decay} = \min \left(\beta, \frac{1 + \text{num_updates}}{10 + \text{num_updates}} \right),$$

where $\beta = 0.9999$ is the asymptotic upper bound. This schedule begins with a low decay (approximately 0.1 at step 0), enabling rapid adaptation of the shadow parameters during early training, and gradually increases toward β as training progresses. The rising decay places increasing emphasis on historical parameters, thereby stabilizing the EMA model in later stages. We use a fixed batch size of 256 and do not apply any data augmentation.

Data preprocessing. We strictly adhere to the data preprocessing protocol of ADM (Dhariwal & Nichol, 2021). Input images are encoded into latent vectors $\mathbf{z} \in \mathbb{R}^{32 \times 32 \times 4}$ using the pretrained VAE from Stable Diffusion (Rombach et al., 2022).

Evaluation details. Following Dhariwal & Nichol (2021), we use the same reference batches and evaluation setup as in the official ADM implementation.¹ For faster sampling, we enable TF32 precision, consistent with the REPA setup.

We evaluate generation quality using the following standard metrics:

¹<https://github.com/openai/guided-diffusion/tree/main/evaluations>

Table 3: Two-stage experiment, which compares the SiT framework and our StableVM framework with bank size 256. Each model is trained for 500k steps, and the checkpoints are evaluated every 40k steps.

# Steps	SiT					StableVM (Bank=256)				
	IS \uparrow	FID \downarrow	sFID \downarrow	Precision \uparrow	Recall \uparrow	IS \uparrow	FID \downarrow	sFID \downarrow	Precision \uparrow	Recall \uparrow
40k (SiT) / 50k (StableVM)	212.5	4.33	6.29	0.71	0.69	219.2	3.92	6.04	0.72	0.69
80k	223.8	3.53	5.61	0.73	0.68	227.3	3.35	5.61	0.73	0.67
120k	234.4	3.07	5.34	0.74	0.67	234.8	3.02	5.40	0.74	0.67
160k	239.7	2.81	5.24	0.74	0.67	238.9	2.87	5.74	0.74	0.66
200k	244.7	2.62	5.20	0.75	0.67	247.0	2.59	5.21	0.75	0.67
240k	247.4	2.54	5.17	0.75	0.67	250.3	2.48	5.20	0.75	0.67
280k	250.9	2.47	5.14	0.75	0.66	253.3	2.40	5.13	0.75	0.67
320k	251.7	2.39	5.17	0.75	0.67	255.4	2.34	5.13	0.75	0.66
360k	255.4	2.35	5.16	0.75	0.67	258.5	2.28	5.11	0.76	0.66
400k	258.6	2.28	5.17	0.76	0.67	261.7	2.17	4.98	0.76	0.67
440k	259.3	2.29	5.20	0.76	0.67	262.3	2.16	5.00	0.76	0.67
480k	259.7	2.26	5.22	0.76	0.67	263.1	2.16	5.03	0.76	0.66
500k	260.1	2.26	5.23	0.76	0.67	262.7	2.17	5.06	0.76	0.67

- **FID** (Heusel et al., 2017): Measures the Fréchet distance between feature distributions of real and generated images using the Inception-v3 network (Szegedy et al., 2016), under the assumption that both are multivariate Gaussian.
- **sFID** (Nash et al., 2021): Computes FID using spatially resolved intermediate features from Inception-v3 to better capture local structure and spatial coherence.
- **IS** (Salimans et al., 2016): Evaluates sample diversity and quality via the Inception-v3 logits, defined as the KL divergence between the marginal label distribution and the conditional distribution of predicted labels.
- **Precision and Recall** (Kynkäänniemi et al., 2019): Quantify sample fidelity (precision) and coverage of the true data manifold (recall) based on nearest-neighbor distances in feature space.

E.1.2 TWO-STAGE TRAINING AND SAMPLING

For the two-stage experiment, we use $\xi = 0.7$ as the split point between the low-variance regime and the high-variance regime. During training, we sample the timestep t uniformly from $[\xi, 1]$ and compute the loss for optimization. We use the CFM loss as the baseline to compare with our StableVM loss. Other configurations is the same as those stated in Appendix E.1.1. In particular, for $t \in [\xi, 1]$, we train an SiT-XL/2 from scratch, and for $t \in [0, \xi]$, we use a pretrained SiT-XL/2 from REPA² for stepping. Table 3 shows the detailed performance results of the two methods. Training is done up to 500k steps, and the checkpoints are evaluated every 40k steps. We note from the results that our model consistently perform better than the SiT with CFM loss.

E.1.3 COMPARISONS WITH BASELINES

We further extend the training duration to 2M iterations and compare our SiT-XL model trained with the CFM and StableVM loss against a range of strong baselines: ADM (Dhariwal & Nichol, 2021), VDM++ (Kingma & Gao, 2024), Simple Diffusion (Hoogeboom et al., 2023), CDM (Ho et al., 2022), LDM (Rombach et al., 2022), U-ViT (Bao et al., 2023), DiffiT (Hatamizadeh et al., 2024), MDTv2 (Gao

²<https://github.com/sihyun-yu/REPA/blob/main/utils.py#14>

Table 4: Comparisons on ImageNet 256×256 with CFG. \downarrow and \uparrow indicate whether lower or higher values are better, respectively.

Model	FID \downarrow	sFID \downarrow	IS \uparrow	Precision \uparrow	Recall \uparrow
<i>Pixel diffusion</i>					
ADM-U	3.94	6.14	186.7	0.82	0.52
VDM++	2.40	-	225.3	-	-
Simple diffusion	2.77	-	211.8	-	-
CDM	4.88	-	158.7	-	-
<i>Latent diffusion, U-Net</i>					
LDM-4	3.60	-	247.7	0.87	0.48
<i>Latent diffusion, Transformer + U-Net hybrid</i>					
U-ViT-H/2	2.29	5.68	263.9	0.82	0.57
DiffiT*	1.73	-	276.5	0.80	0.62
MDTv2-XL/2*	1.58	4.52	314.7	0.79	0.65
<i>Latent diffusion, Transformer</i>					
MaskDiT	2.28	5.67	276.6	0.80	0.61
SD-DiT	3.23	-	-	-	-
DiT-XL/2	2.27	4.60	278.2	0.83	0.57
SiT-XL/2	2.05	4.46	255.4	0.81	0.59
+ StableVM(bank=256)	2.03	4.40	260.1	0.82	0.59

et al., 2023), MaskDiT (Zheng et al., 2024), DiT (Peebles & Xie, 2023), and SiT (Ma et al., 2024). Full results are reported in Table 4. Despite training for only 2M iterations—significantly fewer than the 7M iterations used in Ma et al. (2024)—our model achieves competitive performance and closely matches the results reported for SiT. Moreover, under this training budget, our approach converges faster and attains better performance than the CFM-trained SiT baseline.

E.1.4 UNCONDITIONAL GENERATION

We also evaluate Algorithm 1 in the unconditional generation setting. Specifically, we construct a synthetic Gaussian Mixture Model (GMM) distribution and train the model to learn it using either the standard CFM loss or our proposed StableVM loss.

The GMM is defined with 100 modes. For each component k , the mean vector μ_k is sampled independently from a uniform distribution over $[-1, 1]$ in each dimension. The variances are drawn independently per component and per dimension from a uniform distribution over $[10^{-2}, 10^{-1}]$, yielding anisotropic Gaussian components. The mixing coefficients π are obtained by sampling each entry from $\text{Uniform}(0.1, 1.0)$ and normalizing so that they sum to one. To generate samples, we first draw a component index k via multinomial sampling according to π , then sample from the corresponding Gaussian using the reparameterization trick:

$$\mathbf{x} = \mu_k + \sigma_k \odot \epsilon, \quad \text{where } \epsilon \sim \mathcal{N}(0, \mathbf{I}),$$

and \odot denotes element-wise multiplication.

We fix the data dimensionality to 10 and evaluate both the CFM loss and our proposed StableVM loss with a reference batch size of 2048 on this distribution. Model performance is assessed by computing the *second-order moment* of the discrepancy between the model’s predicted velocity field $v_\theta(\mathbf{x}_t, t)$ and the true velocity

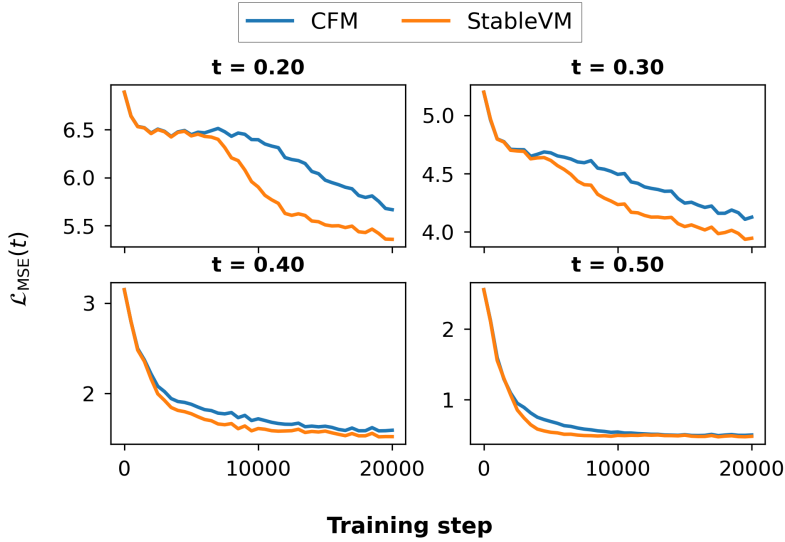


Figure 5: Comparison of StableVM and CFM on the synthetic GMM distribution. We plot the *second-order moment* as a function of training iterations at four time steps: $t = 0.20, 0.30, 0.40$, and 0.50 .

field $\mathbf{v}_t(\mathbf{x}_t)$, under the marginal distribution $p_t(\mathbf{x}_t)$:

$$\mathcal{L}_{\text{MSE}}(t) = \frac{1}{2} \mathbb{E}_{\mathbf{x}_t \sim p_t} \left[\|\mathbf{v}_\theta(\mathbf{x}_t, t) - \mathbf{v}_t(\mathbf{x}_t)\|^2 \right]. \quad (45)$$

However, the exact velocity field $\mathbf{v}_t(\mathbf{x}_t) = \mathbb{E}_{p_t(\mathbf{x}_0|\mathbf{x}_t)}[\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0)]$ is intractable. We therefore approximate it using a self-normalized importance sampling estimator (Hesterberg, 1995). Concretely, given $N = 50,000$ samples $\{\mathbf{x}_0^i\}_{i=1}^N$ drawn from the data distribution, we approximate:

$$\mathbf{v}_t(\mathbf{x}_t) \approx \sum_{i=1}^N w_i(\mathbf{x}_t) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0^i), \quad \text{where} \quad w_i(\mathbf{x}_t) = \frac{p_t(\mathbf{x}_t | \mathbf{x}_0^i)}{\sum_{j=1}^N p_t(\mathbf{x}_t | \mathbf{x}_0^j)}.$$

This approximation enables us to estimate the *second-order moment* at any time step t .

The results are shown in Fig. 5, with separate curves for $t = 0.20, 0.30, 0.40$, and 0.50 . The plots clearly demonstrate that StableVM improves both the convergence speed and the final performance compared to standard CFM.

Reproduction of Fig. 1. To reproduce the variance curves in Fig. 1, we follow the same procedure described above but increase the dimensionality of the synthetic GMM to 100 and 500, using identical configuration settings. The true velocity field is again approximated using the self-normalized importance estimator with $N = 10,000$ samples. For the CIFAR-10 curve, we use the full training set of 50,000 samples for evaluation.

Table 5: Evaluation results of Wan2.2 on video generation with 25 steps. Higher is better (\uparrow) except for LPIPS (\downarrow).

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
25 steps(Ours)	28.874	0.907	0.066
25 steps	12.379	0.477	0.551

E.2 STABLE VELOCITY SAMPLING

E.2.1 EXPERIMENTAL SETUP

We implement our algorithm primarily using the Huggingface `diffusers` library (von Platen et al.), which supports state-of-the-art pretrained models. Following the recommended configurations, we adopt guidance scales of 3.5, 5.0, and 4.5 for Flux (Labs, 2024), SD3 (Esser et al., 2024), and SD3.5 (Esser et al., 2024), respectively. For Wan2.2 (Wan et al., 2025), we set the `guidance scale` to 4.0 and `guidance scale 2` to 3.0.

Our algorithm introduces three key hyperparameters: the split point ξ , the number of steps in the low-variance regime NFE_{Low} , and the variance factor f_β . As discussed in Sec. 3.2 and Eq. (12): (i) the split point ξ determines the boundary between the *low-variance* and *high-variance* regimes, (ii) the variance factor f_β controls the variance of the DDIM-style posterior in Eq. (12), and (iii) NFE_{Low} specifies the number of sampling steps allocated to the *low-variance regime*. In practice, we partition the default time scheduler of pretrained models into two regimes according to ξ . The low-variance regime is then uniformly divided into NFE_{Low} steps, and sampling is performed following Eq. (12). For Tab. 2, we set $\xi = 0.85$, $\text{NFE}_{\text{Low}} = 8$, and $f_\beta = 0$.

Evaluation Details. For text-to-image tasks, we evaluate different samplers on the *Geneval* benchmark (Ghosh et al., 2023), which consists of 553 prompts spanning six categories. For each experiment, we generate 4 samples using random seeds $\{0, 1000, 2000, 3000\}$, strictly following the official *Geneval* evaluation protocol.

For text-to-video tasks, we employ Wan2.2 to generate videos conditioned on 100 prompts from the `human_action` dimension of VBench (Huang et al., 2024). We report only reference-based metrics by comparing generated videos against 50-step baseline results. All generations are performed with the same random seed for consistency. Table 5 presents the evaluation results of Wan2.2 on video generation with 25 steps. Our custom sampler achieves significantly higher PSNR and SSIM and much lower LPIPS, indicating that it produces results much closer to the 50-step baseline compared to the original sampler.

E.2.2 ABLATION STUDIES

Table 6 presents the effect of three key hyperparameters—variance factor f_β , the number of low-variance steps NFE_{Low} , and the split point ξ —on Stable Velocity Sampling with SD3.5 at 1024×1024 .

Compared to the baseline, introducing a split at $\xi = 0.85$ with $\text{NFE}_{\text{Low}} = 8$ already leads to large gains across all metrics (PSNR from 15.58 to 30.95, LPIPS from 0.370 to 0.050), showing the benefit of separating low- and high-variance regimes.

Varying the variance factor f_β reveals a trade-off: increasing f_β slightly reduces PSNR but improves overall stability, with $f_\beta = 0.2$ achieving the best overall score (0.722).

Table 6: Ablation study on split point ξ , low-variance regime steps NFE_{Low} , and variance factor f_β of Stable Velocity Sampling for SD3.5 at 1024×1024 . Higher is better (\uparrow), except for LPIPS (\downarrow).

NFE	ξ	NFE_{Low}	f_β	Overall \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Baseline							
25	–	–	–	0.711	15.576	0.735	0.370
25	0.85	8	0.0	0.718	30.954	0.958	0.050
Variance factor f_β							
25	0.85	8	0.1	0.719	30.094	0.945	0.056
25	0.85	8	0.2	0.722	28.536	0.915	0.071
Low-variance regime steps NFE_{Low}							
20	0.85	3	0.0	0.711	26.220	0.873	0.155
30	0.85	13	0.0	0.719	33.891	0.977	0.025
Split point ξ							
20	0.90	8	0.0	0.727	26.171	0.918	0.098
29	0.80	8	0.0	0.721	34.387	0.974	0.031
35	0.70	8	0.0	0.721	39.303	0.987	0.014

For the number of low-variance steps, allocating more steps consistently improves fidelity: with $\text{NFE}_{\text{Low}} = 13$, PSNR rises to 33.89 and LPIPS drops to 0.025.

Finally, the split point ξ plays a critical role. An earlier split at $\xi = 0.70$ with 35 steps achieves the strongest results overall, reaching the highest PSNR (39.30), SSIM (0.987), and lowest LPIPS (0.014). This confirms that ξ enables the sampler to fully exploit low-variance dynamics for both fidelity and perceptual quality.

F LLM USAGE STATEMENT

In preparing this manuscript, we used OpenAI’s ChatGPT solely as a language editing tool to polish grammar, phrasing, and clarity of our writing. The model was not involved in research ideation, experimental design, analysis, or drawing conclusions. All scientific content and interpretations are the authors’ own.

G MORE QUALITATIVE RESULTS

We provide additional qualitative results for SD3.5 (Esser et al., 2024), Flux (Labs, 2024), SD3 (Esser et al., 2024), and Wan2.2 (Wan et al., 2025), shown in Fig. 6, Fig. 7, Fig. 8, and Fig. 9, respectively.

1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597

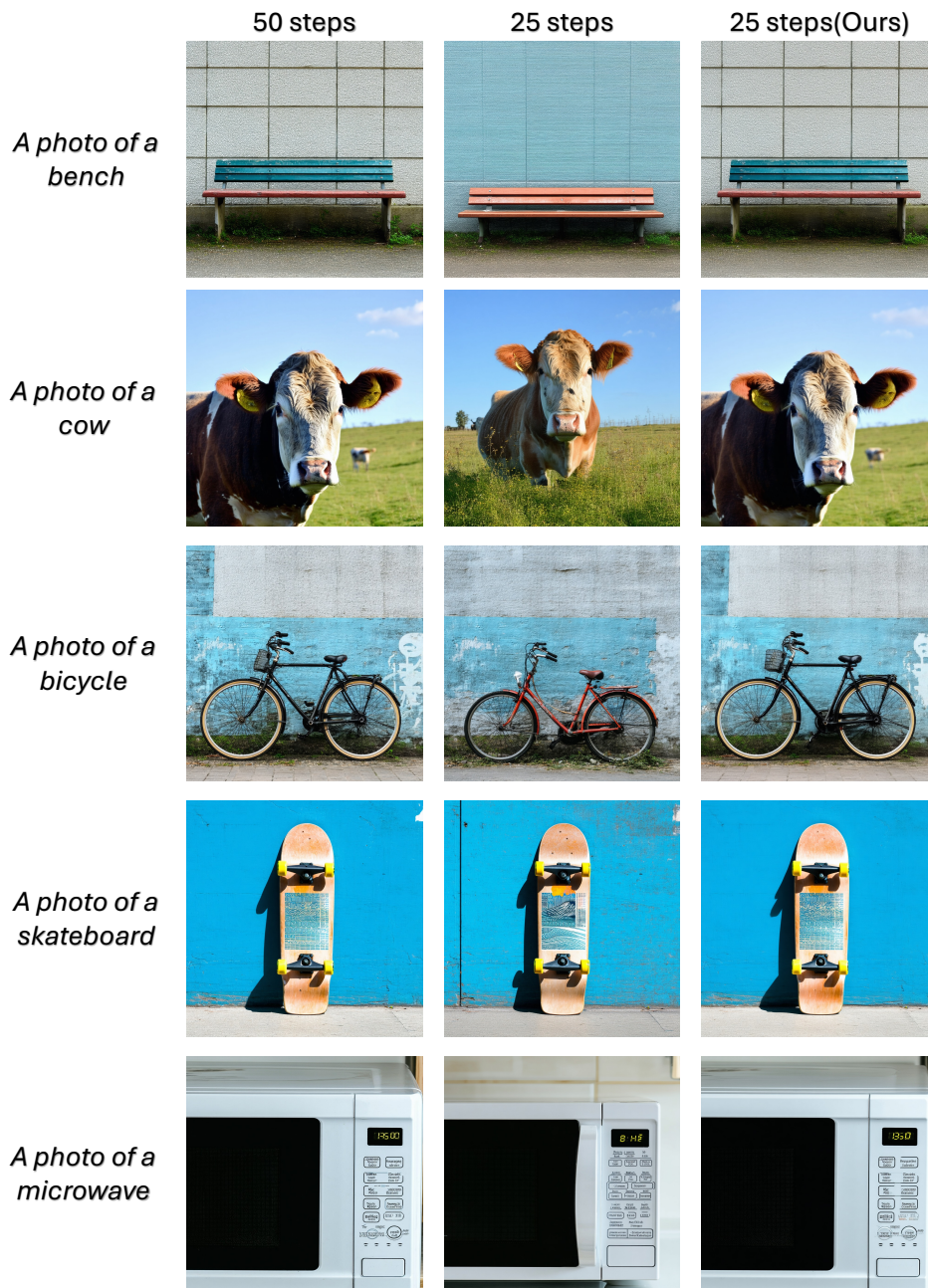


Figure 6: Visual comparison of SD3.5-Large (Esser et al., 2024) on different prompts. We show results using the default sampler with 50 and 25 steps, alongside our 25-step sampler under the same random seed. Zoom in for finer details.

1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644

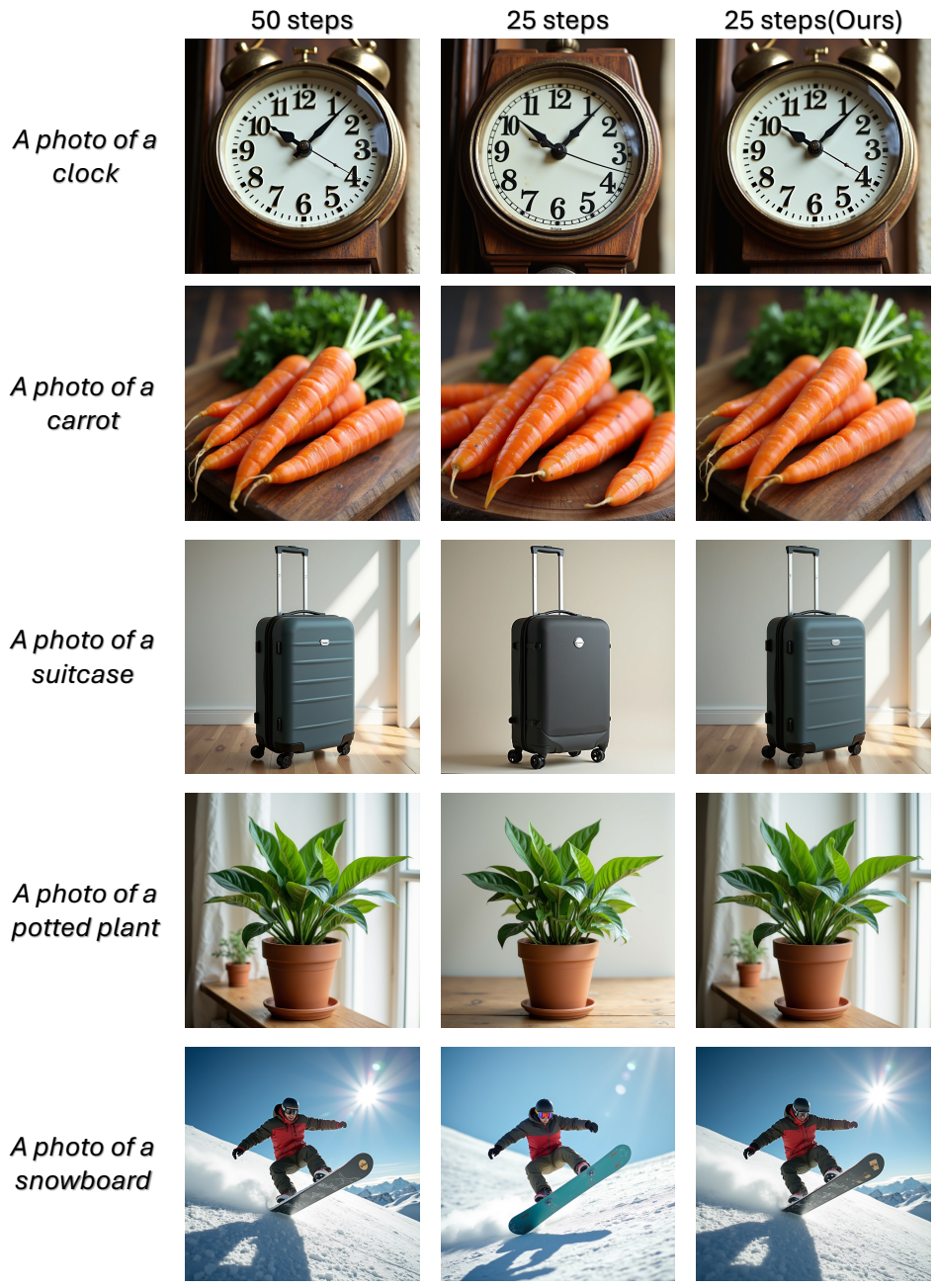


Figure 7: Visual comparison of Flux-dev (Labs, 2024) on different prompts. We show results using the default sampler with 50 and 25 steps, alongside our 25-step sampler under the same random seed. Zoom in for finer details.

1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673
 1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691



Figure 8: Visual comparison of SD3-medium (Esser et al., 2024) on different prompts. We show results using the default sampler with 50 and 25 steps, alongside our 25-step sampler under the same random seed. Zoom in for finer details.

1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738

Prompt: A person is air drumming.

50 steps



25 steps



25 steps
(Ours)



Prompt: A person is brushing teeth.

50 steps



25 steps



25 steps
(Ours)



Figure 9: Visual comparison of Wan2.2 (Wan et al., 2025) on different prompts. We show results using the default sampler with 50 and 25 steps, alongside our 25-step sampler under the same random seed. Zoom in for finer details.